

AD-A054 556

WISCONSIN UNIV-MADISON MATHEMATICS RESEARCH CENTER
STUDIES IN THE ANALYSIS OF SERIALY DEPENDENT DATA.(U)
MAR 78 L C PALLESEN

F/G 12/1

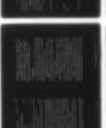
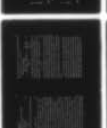
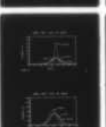
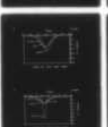
DAA629-75-C-0024

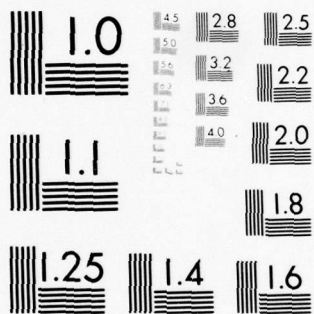
UNCLASSIFIED

MRC-TSR-1837

NL

1 OF 2
AD
A054556



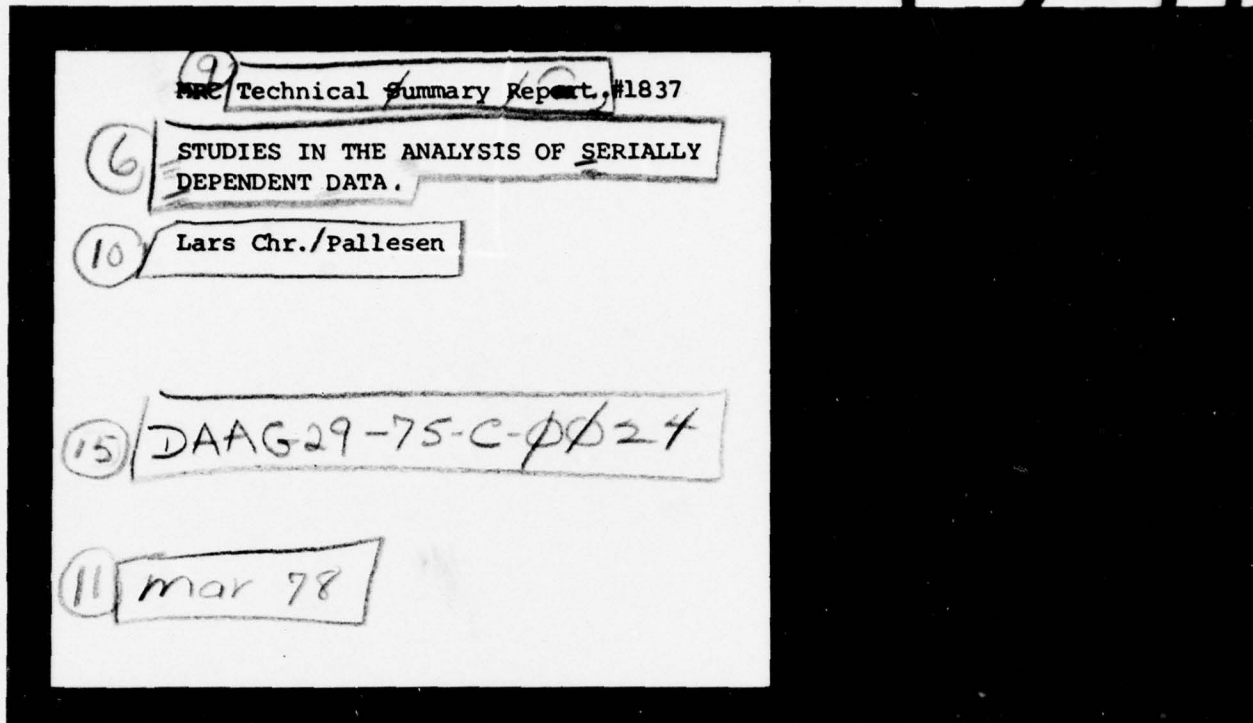


MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

FOR FURTHER TRAN *FILE*

(13) *R*

AD A 054556



Mathematics Research Center
University of Wisconsin-Madison
610 Walnut Street
Madison, Wisconsin 53706

(14) MRC-TSR-1837

March 1978

(12) 97p.

Received July 8, 1977

DDC FILE COPY

DDC
REF ID:
JUN 2 1978
D

Approved for public release
Distribution unlimited

Sponsored by

U. S. Army Research Office
P.O. Box 12211
Research Triangle Park
North Carolina 27709

221 200

Gu

UNIVERSITY OF WISCONSIN-MADISON
MATHEMATICS RESEARCH CENTER

STUDIES IN THE ANALYSIS OF SERIALLY DEPENDENT DATA

Lars Chr. Pallesen

Technical Summary Report #1837
March 1978

ABSTRACT

The analysis of linear models with independent homoscedastic, normal noise (white noise) occupies a prominent position in applied statistics. This report is concerned with the linear model analysis of data which cannot be assumed statistically independent because the data have been collected sequentially in time or space.

Assuming that the noise in a linear model (with p -dimensional parameter vector $\hat{\theta}$) follows a first order autoregressive scheme (with parameter ϕ) it is shown in Chapter 2 how in the Bayesian framework inferences can be drawn about $\hat{\theta}$ and ϕ jointly, conditionally and marginally. Two AR-1 schemes are considered, one covering explosive as well as stationary situations, the other assumes stationarity and reversibility a priori. The important task of choosing an appropriate joint prior distribution for the parameters is given special attention.

Recognizing that the assumption of independence can be a crucial one, it has become a widespread practice in much regression work, where observations are made in time order, to test for serial correlation using the well known Durbin-Watson test. This test corresponds to determining whether a certain estimate $\hat{\phi}$ of ϕ is significantly different from zero. It is shown in Chapter 3 that only in relation to a model lacking a mean do suggested alternative testing procedures show decisively greater empirical power than the DW-test. However it is argued that tests of a null hypothesis of independence should really not be carried out when serial correlation is to be expected. It is demonstrated that inferences about $\hat{\theta}$ may be quite misleading if made conditionally on $\hat{\phi} = 0$.

Sponsored by the United States Army under Contract No. DANG29-75-C-0024.

(when that hypothesis is "accepted") or conditionally on $\hat{\phi}$ (when it is "rejected"). The Bayesian analysis does not suffer from these handicaps. The Bayesian inference about $\hat{\theta}$ may be approximated by a conditional inference using the maximum likelihood estimate $\hat{\phi}$ of ϕ ($\hat{\phi}$ is found by estimating $\hat{\phi}$ and $\hat{\theta}$ simultaneously by least squares), and this conditional analysis is paralleled in sampling theory framework when inferences about $\hat{\theta}$ are drawn as if $\hat{\phi} = \hat{\phi}$.

One way in which data analysts have sometimes tried to get around the problem of dependence, is by differencing and then assuming that the errors of the differences are independent. This is equivalent to assuming that $\phi = 1$ in the AR-1 noise model; and the plausibility of this particular value for ϕ may be assessed by studying the marginal posterior distribution of ϕ . For models involving a mean this parameter vanishes for $\phi = 1$. This creates a singularity, and it is shown in Chapter 4, how the density may be found at such (distinct) points. It is also developed in the Bayesian framework, how the appropriateness of differencing may be assessed using posterior model probabilities for two alternative noise models, one assuming the original observations to have AR-1-noise, the other that the differences have AR-1 noise.

AMS (MOS) Subject Classification: Primary 62F15, 62M10
Secondary 62P20, 62N99, 62J05

Key words: Regression, Linear models, Bayesian approach,
Serial correlation, Autoregressive noise,
Nonstationary noise, Differencing

Work Unit No. 4 - Probability, Statistics and Combinatorics

ACCESSION NO.	
NTIS	White Section <input checked="" type="checkbox"/>
DOC	Grey Section <input type="checkbox"/>
UNANNOUNCED	<input type="checkbox"/>
JUSTIFICATION	
BY	
DISTRIBUTION/AVAILABILITY CODES	
Ext.	AVAIL. and/or SPECIAL
A	

SIGNIFICANCE AND EXPLANATION

One of the most widely used statistical techniques in applied fields like engineering, business, economics, sociology etc. is linear model analysis (analysis of variance, regression analysis, etc.). When the data to be analyzed are collected sequentially in time or space, it turns out, more often than not, that the data at any given point in space or time depend to some extent on the data at neighboring points in space or in instants of time. The data are then said to be serially correlated (autocorrelated). However "standard" methods of analysis usually assume the observations to be statistically independent; if the assumption of independence is incorrect, this will invalidate the analysis and may cause very misleading conclusions.

This report is devoted to studying how serially dependent data may be analyzed, when one extra parameter is incorporated to allow for serial correlation. Specifically it is developed in the Bayesian framework how general linear models may be analyzed when the noise follows a first order autoregressive process (a Markov process). We argue against following the popular practice of checking independence with a Durbin-Watson test whenever serial correlation is feared.

One way to get around the problems of serial correlation which has sometimes appeared helpful is to difference the series, i.e. analyze the increments rather than the original data. The question of whether to difference or not to difference is looked into in a quantitative way.

ACKNOWLEDGEMENTS

Professor G. E. P. Box served as my advisor during the preparation of this thesis, and I deeply appreciate his generous sharing of ideas, humor, wisdom and dedication. I am proud and grateful to count myself among his disciples.

I should also like to express gratitude towards my diverse financial sponsors over the years; in particular I shall acknowledge the vital grants from "Otto Mønstedts Fond" and the never failing backing by "Civilingeniør Frants Allings Legat".

The support in many ways of our close family and our dear friends has been invaluable.

I thank my wife for making my labor worthwhile.

CONTENTS

1. Introduction	1	Appendix B, Consequence of assuming prior independence between θ and ϕ	71
1.1 Models and parametric inference	1	3. Probing for Serial Correlation in Least Squares Regression.	73
1.2 Broadening the class of linear models	3	3.1 The Durbin-Watson test	75
1.3 Analysis of serially correlated data.	4	3.2 $SS(\hat{\theta}, \hat{\phi})$ as a function of ϕ	78
1.4 Outline of present work	7	3.2.1 Singularity due to a mean	78
2. Linear Models with AR-1 Noise: A Bayesian Analysis.	11	3.2.2 Simple singularity.	80
2.1 The class of models	12	3.2.3 General singularity	81
2.2 Likelihood functions.	16	3.2.4 Singularity in polynomial regression.	83
2.2.1 Derivation of likelihood functions	16	3.2.5 Summary of findings	84
2.2.2 Maximum likelihood estimates	18	3.3 Inference about ϕ based on maximized likelihood.	85
2.2.3 Fisher's information matrix.	21	3.4 Bayesian inference about ϕ	88
2.3 Prior distributions	22	3.5 Examples	91
2.3.1 Prior for the a priori stationary model.	23	3.5.1 The data generated by Zellner and Tiao.	91
2.3.2 Prior for the not necessarily stationary model	27	3.5.2 The textile data.	99
2.4 Posterior distributions	28	3.5.3 The spirits data.	104
2.4.1 Marginal posterior distribution for (θ, ϕ) jointly.	29	3.6 Monte Carlo study.	110
2.4.2 Marginal posterior distribution of ϕ	29	3.6.1 Power curves.	112
2.4.3 Marginal posterior distribution of θ	30	3.6.2 Size curves	119
2.5 Marginal prior and posterior for ϕ in the not necessarily stationary model.	32	3.7 Conclusion	123
2.6 Examples.	41	Appendix. Additional results from the Monte Carlo study.	125
2.6.1 An artificial example.	41	4. To Difference or not to Difference Data Series in Linear Model Analysis: A Bayesian Study	132
2.6.2 Detecting a change in level.	52	4.1 Joint posterior distribution of (θ_j, ϕ_0)	134
2.7 Conclusion.	57	4.2 Posterior density of the AR-1 parameter at unity	137
Appendix A, On the noninformative prior for ϕ	60		

4.3	Expending one degree of freedom on differencing rather than on a fixed mean?	141
4.3.1	The alternative likelihood function	141
4.3.2	The complementing priors	144
4.3.3	Posterior model probabilities	147
4.4	The CGK data	149
4.5	Conclusion	171
	Appendix. Other applications	173
5.	Summary	175
	References	179

CHAPTER 1

Introduction.

This introductory chapter contains a general discussion of some issues relating to the topics of this thesis, viz: Inference about the linear parameters θ_j , $j = 1, 2, \dots, p$, of a general linear model with first order autoregressive (AR-1) noise, inference about its autoregressive parameter ϕ , inference about $\{\theta_j, \phi\}$ jointly, and quantitative assessment of the appropriateness of differencing data series prior to such a linear model analysis.

1.1. Models and parametric inference.

There are two requirements a statistical model should satisfy in order to be useful in applications. First it must sufficiently approximate reality; second it should be parsimonious, that is it must be simple enough to lend itself to analysis and interpretation. If this latter requirement is not satisfied, even a realistic model may not be particularly helpful. When such a model is exceedingly complex involving a vast number of parameters, (i) it may be impossible to verify, (ii) it may contain many more unknown parameters than can be estimated with available data, and (iii) its implications may be beyond comprehension anyway.

Of course there are limits to how simple a model can be and still be useful, so it may not always be possible to identify an adequate model among the "standard" model forms, for which satisfactory techniques for analysis have been developed. The purpose of the present work, is to contribute to extending the "stock" of standard model structures, which can be handled satisfactory from a statistical point of view.

The usefulness of a parametric model in applications may be attributed to its ability to extract the relevant information embodied in a data set, and make it accessible for human understanding and use. Typically a useful statistical model incorporates more parameters than those of primary interest. These additional parameters serve in making the model a realistic one, and are sometimes called "nuisance parameters". These nuisance parameters are needed if the inference about the primary parameters is to be valid, but their inclusion has the side effect of complicating the analysis of the model. In the Bayesian framework, (see for example Jeffreys [1961], Box and Tiao [1973]), the difficulties created by nuisance parameters are mainly of a computational nature, since inferences about any subset of parameters can be made from their marginal posterior distribution. This distribution is found from the joint posterior by integrating out those parameters not under consideration. Frequently this integration cannot be carried out analytically, but in principle it can always be done numerically using an electronic computer.

Throughout the present work a Bayesian approach is taken almost exclusively. In this framework the choice of an appropriate prior distribution for the parameters is part of the model specification. The prior may serve as a carrier of already existing information or to incorporate prior judgement in the analysis. However in scientific investigations, one would usually want the prior to be "neutral" or "noninformative" relative to the information in the data. It is not always clear which form the prior should have to express such a desire, and this is among the questions we shall address in this thesis.

appropriate after some suitable power transformation has been applied to the dependent variable. They demonstrated how inference can be made about the transformation and in turn about the linear parameters $\underline{\theta}$, by studying the maximized likelihood function as a function of the transformation parameter(s) λ ($\underline{\lambda}$) or by computing the marginal posterior distribution of λ ($\underline{\lambda}$).

In some situation the linear model (1.2) should realistically be generalized to allow for the possibility, that one or more of the observations are outliers, i.e. do not conform to the model of the remaining "good" observations. Again the literature on outliers is extensive (in particular see Anscombe [1960], Anscombe and Tukey [1963] and Tukey [1960]). From the Bayesian viewpoint Box and Tiao [1968] assumed in their analysis that a "bad" observation has the same expected value as a good one, but with an inflated variance. Abraham and Box [1975] assumed in an alternative Bayesian analysis of the same problem, that a shift in mean with no change in variance, causes a "bad" value.

It is the purpose of the present work to address some questions relating to the analysis of serially correlated data, i.e. when the augmentation of the linear model should be of the third kind mentioned above.

1.3 Analysis of serially correlated data.

Whenever data are collected sequentially in space or in time, serial correlation is to be expected because of the natural continuity of matter and of events developing in time. Thus it is of interest to broaden the class of linear models (1.2) by relaxing the assumption of independence.

At one time there was rather little concern about such dependencies, in particular regression analysis of time series data were commonly conducted under the assumption, that the observations had white noise

1.2 Broadening the class of linear models.

Writing a model for a dependent variable y_i , $i = 1, 2, \dots, n$, as

$$\underline{y} = \underline{X} \underline{\theta} + \underline{\varepsilon} \quad (1.1)$$

the statistical analysis is particularly simple if the deterministic part of the model \underline{X} is linear in the parameters, i.e. $\underline{y} = \underline{X} \underline{\theta}$, where \underline{X} is the matrix of independent variables and where $\underline{\theta}$ is the p -dimensional vector of linear parameters; and if it is assumed that the stochastic part of the model $\underline{\varepsilon}$ is equal to $\underline{\varepsilon}$, where the ε_i 's, $i = 1, 2, \dots, n$, are i.i.d. $N(0, \sigma^2)$ (white noise):

$$\underline{y} = \underline{X} \underline{\theta} + \underline{\varepsilon} \quad (1.2)$$

Frequently the assumptions about \underline{y} implied by this linear model are found too restrictive to be realistic for data sets met in practice.

An adequate representation of the data at hand may often be achieved however, if the class of linear models is broadened by parsimonious use of additional parameters.

To do this we must consider which are most likely ways for the model to be wrong. Box [1976] has suggested that three prime contenders are

- (i) need for transformation of the dependent variable,
- (ii) allowance for outliers,
- (iii) allowance for serial correlation when data are collected sequentially.

The necessary modification in the analysis of these augmentations of the linear model are readily derived using a Bayesian approach.

The idea of transformation of original data is old (see in particular Bartlett [1947] and Tukey [1957]). From the Bayesian viewpoint Box and Cox [1964] analyzed the situation when the model form (1.2) is

errors. Surveys of the literature relating to the effects of serial correlation in regression are given by Anderson [1954] and by Watson [1967]. Realizing that serial correlation could be important Durbin and Watson [1950], [1951] produced a test (the DW-test) which over the years has become widely used. This test is in effect equivalent to estimating a first order autoregressive (AR-1) parameter $\hat{\phi}$ from the residuals left by an ordinary least squares fit, and seeing whether this estimate $\hat{\phi}$ is significantly different from zero. Modifications and new tests have been proposed for example by Abrahamse and Louter [1971], Berenblut and Webb [1973], Durbin [1969], [1970], Gray [1970] and Schmidt [1972], but from the point of view of power, none of these alternatives have been demonstrated to have any decisive advantage over the DW-test, Durbin and Watson [1970], Smith [1976]. Durbin [1970a] argues, that the main incentive for testing a null hypothesis of independence is that the analysis becomes simple if that hypothesis can be "accepted". When it is "rejected", Theil and Nagar [1961] suggest to make inference about the linear parameters $\hat{\beta}$ assuming the noise model to be AR-1 with $\hat{\phi} = \hat{\phi}$ (calculated from the DW statistic).

Cochrane and Orcutt [1949] demonstrated, that ordinary least squares estimates of regression parameters may be very poor if serial correlation is present, and that the situation may be remedied by applying an independence inducing transformation. If the autocorrelation structure is not known, they express scepticism (see also Orcutt and Cochrane [1949]) about the possibility of estimating the structural parameters from ordinary least squares residuals, since they are biased towards randomness: but an iterative procedure for estimating one or two autoregressive parameters is outlined.

Chapmanowne [1948] suggested a Bayesian approach to estimating the parameters of a p -th order autoregressive noise process in a regression model; and Durbin [1960] proposed an asymptotically efficient two step estimation procedure of the regression parameters in a similar model, taking the autoregressive parameters into account. For regression models with AR-1 noise Hildreth [1969] showed, that the maximum likelihood estimates of the parameters are asymptotically efficient and equal to those of Durbin [1960].

Reiser [1975] analyzed the regression model with p -th order autoregressive noise applying a structural inference approach, see Fraser [1968]. In this framework inference about (including estimation of) the autoregressive parameters $\hat{\phi}$ is made using a marginal likelihood function, and inferences about the regression parameters $\hat{\beta}$ are drawn conditional on the estimates of $\hat{\phi}$. The structural inference approach to estimation of autoregressive parameters has also been taken by Haq [1970], [1971] and Levenbach [1972].

Zellner and Tiao [1964] showed how marginal inference can be made about regression parameters when the noise follows an AR-1 scheme with unknown parameter $\hat{\phi}$. Sredni [1970] generalized this Bayesian analysis to cover a p -th order autoregressive noise process.

The noise N_1 in (1.1) might be considered generated by a member of the very wide class of autoregressive integrated moving average (ARIMA) processes proposed by Box and Jenkins [1970]. Their work suggests that an appropriate ARIMA model of order (p,d,q) should be specified for N_1 , and then the combined model consisting of both a deterministic (linear) part and a stochastic part should be analyzed. Differencing the data series d times produces a new series whose errors follow an

ARMA(p,q) model.

Pierce [1971] showed, that estimating the regression parameters and the stochastic parameters of an ARMA(p,q) noise model simultaneously by least squares asymptotically produces maximum likelihood estimates (assuming Normality). The estimates of the regression parameters are asymptotically uncorrelated with the ARMA parameter estimates, and the latter have the same covariance matrix as would have resulted if the deterministic model component did not exist.

The qualitative use of sample autocorrelation functions to determine when differencing is appropriate, is discussed by Box and Jenkins [1970]. It was pointed out very early by Student [1914], that differencing of sequential data could be employed to eliminate spurious cross correlation between series of observations. Orcutt and Cochrane [1949] suggested, that it might often be helpful in the analysis of economic time series to look at the first differences rather than the original data series. Recently Anderson [1975] studied when the variation of low order ARMA (p,q) processes is reduced by differencing.

1.4 Outline of present work.

As far as the present thesis is concerned, attention is limited to situations where the autocorrelated errors may adequately be represented by a one parameter noise model as specified below.

A naturally appealing noise model would be one where the serial correlation fell off smoothly with the distance (or lag) between observations, and a reasonable way for this to happen would seem to be exponentially. Exactly this behavior is obtained from the first order Markov (AR-1) model

$$e_i = \phi e_{i-1} + \epsilon_i \quad (1.3)$$

when $0 < \phi < 1$ (as before the ϵ_i 's are a white noise series). For this model the theoretical autocorrelations are $\rho_1 = \phi$, $\rho_2 = \phi^2$, ..., $\rho_k = \phi^k$, ...

Adopting the noise model (1.3) we have

$$\tilde{y} = X\tilde{\theta} + e \quad (1.4)$$

and of course this less restrictive linear model degenerates to the usual independent model (1.2) when the autoregressive parameter ϕ is zero. If ϕ was known a priori, then the corresponding independence inducing transformation could at once be applied, and inference about the linear parameters $\tilde{\theta}$ could be made by traditional means, Aitken [1935]; but such prior knowledge would hardly ever exist for real data sets.

Chapter 2 contains a Bayesian analysis of the linear model (1.4). This analysis is a continuation of the work of Zellner and Tiao [1964] and Srađni [1970], and it specifically looks into some difficulties which have arisen about choosing a prior distribution for the parameters. Two AR-1 schemes are considered, one covers explosive as well as stationary situations, the other assumes stationarity and reversibility a priori. It is developed how inferences can be drawn about $\tilde{\theta}$ and ϕ jointly, conditionally and marginally.

Chapter 3 examines alternative ways of probing for serial correlation, viz. the DW testing approach, a likelihood approach and a Bayesian approach. A brief inquiry into the relative power of these measures indicates that in this respect the DW-test is quite good. However our study suggests that the DW testing approach as usually carried out has serious shortcomings of quite a different character. It is demonstrated, (1) that correlation sufficiently large to be serious may not be

detected, and that when this is so inference about $\hat{\theta}$ conditional on an "accepted" hypothesis of $\hat{\phi} = 0$ may give very misleading results.

(ii) That, if $\hat{\phi} = 0$ is "rejected" and inference about $\hat{\theta}$ is based on $\hat{\phi} = \hat{\phi}$ (as suggested for example by Thiel and Nagar [1961]) then erroneous conclusions may also be drawn, because $\hat{\phi}$ can be a very poor estimate of ϕ . (iii) That the likelihood and the Bayesian approaches produce not only better estimates for $\hat{\phi}$, but also approximate confidence or HPD intervals for $\hat{\phi}$. (iv) That testing should not be carried out under circumstances where serial correlation is to be expected a priori, but inferences should be made from the less restrictive model (1.4). (v) And that the Bayesian inference about the linear parameters $\hat{\theta}$ may be approximated by a conditional Bayesian analysis using the maximum likelihood estimate $\hat{\phi}$ of ϕ ; and this conditional analysis is paralleled from a sampling theory point of view if inferences about $\hat{\theta}$ are drawn as if $\hat{\phi} = \hat{\phi}$.

The reason why the DW procedure is inappropriate is further clarified by the Bayesian analysis in Chapter 4 of a particular data set, where the joint posterior distributions of (θ_j, ϕ) $j = 1, 2, \dots, p$ are studied.

One way in which data analysts classically have tried to get around the problems of serial correlation is by analyzing the first differences of the data rather than the original data series themselves; and in some instances this practice has appeared to be helpful. This raises the question of whether to difference or not to difference autocorrelated data series, and this is the main topic of Chapter 4. The AR-1 noise model (1.3) reduces to a random walk model as a special case for $\phi = 1$. So if indeed $\phi = 1$ then the increments $\nabla y_t = y_t - y_{t-1}$

have white noise errors and differencing is the proper course of action. Hence it is of interest to determine how plausible that particular value for $\hat{\phi}$ appears in light of the data, and this is looked into from a Bayesian point of view. Of course the increments, ∇y_t , may themselves be autocorrelated, and it may be wondered to what extent the data are more consistent with a model having a fixed mean and an AR-1 noise, or with a model lacking a mean and which implies that the first differences have AR-1 noise. A Bayesian assessment of the relative support of these two alternative noise models is also supplied, which may help decide in which way a single parameter for autocorrelation may best be employed.

CHAPTER 2

Linear Models with AR-1 Noise: A Bayesian Analysis.

In this chapter we conduct a Bayesian analysis of general linear models, when the noise follows a first order autoregressive scheme. The scope of the analysis is that of parametric inference, primarily about the linear parameters, but also about the autoregressive parameter, although this topic will be treated further in Chapter 3.

This problem was studied earlier by Zellner and Tiao [1964]. They demonstrated how marginal inference about linear regression parameters is possible in the Bayesian framework. Their noise model covering explosive cases as well as stationary ones is nearly equivalent to the one analyzed in the following as far as the likelihood is concerned. However we shall argue, that the prior they employed leads to unacceptable consequences which we have remedied.

Szedni [1970] generalized the noise model to a p -th order autoregressive process, but constrained the parameters a priori to cover only stationary situations. On this condition he established the exact likelihood function and an asymptotic noninformative prior for the autoregressive parameters. He did however assume prior independence between the autoregressive parameters on the one hand and the regression parameters on the other hand. Again we believe that this combined prior is inappropriate.

In this chapter the Bayesian treatment of the problem is carried on, with special attention paid to the important task of choosing an appropriate prior distribution. The class of models being considered is a very wide one indeed, covering general linear models, with AR-1 noise which may be stationary or explosive. A subclass (throughout

identified by superscript **) resting upon an a priori assumption of stationarity is treated separately. Two examples are given for illustration.

2.1 The class of models.

We consider the representation of a data vector \underline{y} of n observations by the class of models

$$\underline{y} = \underline{\eta} + \underline{N} \quad (2.1)$$

$n \times 1 \quad n \times 1 \quad n \times 1$

where $\underline{\eta}$ is the deterministic part of the model and \underline{N} is the stochastic part of the model. Of course \underline{y} need not necessarily be the "original" observations, but some transformation of them for example a series of logarithms, or perhaps of first differences.

The structure of $\underline{\eta}$ is supposed specified by the general linear form

$$\underline{\eta} = \underline{X} \underline{\theta} \quad (2.2)$$

$n \times 1 \quad n \times p \quad p \times 1$

where $\underline{\theta}' = [\theta_1 \theta_2 \dots \theta_p]$ is the vector of p linear parameters, and where

$$\underline{X} = [\underline{x}_1 \dots \underline{x}_p] = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix} \quad (2.3)$$

is the matrix of the p independent variables which are supposed known exactly. Like the dependent variable, they need not be the "original" variables, but each may be any given function of the "original" ones; and/or they may be indicator variables. Evidently analysis of variance, analysis of covariance, and regression models are all special cases of this general linear model. Furthermore without loss of generality we

can assume that X has full column rank p .

Concerning the stochastic part of (2.1), we shall consider the

case

$$\tilde{N} = \frac{e}{n \times 1} \quad (2.4)$$

where the random errors, e_i , follow a first order autoregressive scheme

$$e_i = \phi e_{i-1} + \epsilon_i \quad i = 2, 3, \dots, n \quad (2.5)$$

where it is assumed that the ϵ_i 's are i.i.d. $N(0, \sigma^2)$.

The first order autoregressive model has mostly been studied over the stationary range $|\phi| < 1$. If stationarity can be assumed a priori, and if further (see for example Anderson (1954)) the error of the first observation is set

$$e_1 = (1 - \phi^2)^{-1/2} \epsilon_1 \quad (2.6)$$

so that

$$\tilde{e} = \frac{a}{n} e \quad (2.7)$$

where

$$a = \frac{1}{n \times n} \begin{bmatrix} \sqrt{1-\phi^2} & & & 0 \\ -\phi & 1 & & 0 \\ & & \ddots & \\ 0 & & & 1 \end{bmatrix} \quad (2.8)$$

then it is seen that e has a multivariate normal distribution $N(0, I)$

with

$$\tilde{I}^0 = \sigma^2 a^{-1} I a^{-1} = \frac{\sigma^2}{1-\phi^2} \begin{bmatrix} 1 & & & 0 \\ \phi & 1 & & 0 \\ & & \ddots & \\ \phi^{n-1} & & & 1 \end{bmatrix} \stackrel{\text{sym}}{=} \sigma^2 A \quad (2.9)$$

where

$$A = a' a. \quad (2.10)$$

The doubly symmetric feature of \tilde{I}^0 (2.9) shows, that (2.6) implies the error series e_i , $i = 1, 2, \dots, n$, to be reversible, that is the series taken in reverse order is an identical AR-1 process.

The idea of stationarity, i.e. ϕ being between -1 and +1 is really a mathematical concept, but as far as data is concerned there is no particular reason for a cutoff at ± 1 ; and one should be able to deal with explosive situations where ϕ is in fact larger than one. This aspect of the AR-1 model seems however to have been largely ignored.

If stationarity cannot be assumed a priori, instead of (2.5) we write:

$$e_1 = M + \epsilon_1 \quad (2.11)$$

where M is a starting parameter. In this formulation, the autoregressive parameter ϕ may take values on the entire real line. If the process is explosive, then M may be viewed as a measure of how far the explosion has progressed when the first observation is taken. If the process is stationary, then M allows for the possibility that the first error, e_1 , perhaps has an atypical size.

Combining (2.5) and (2.11) we may write

$$\tilde{e} = \frac{b}{n} e - \begin{bmatrix} M \\ 0 \end{bmatrix} \quad (2.12)$$

where

$$b = \frac{1}{n \times n} \begin{bmatrix} 1 & & & 0 \\ -\phi & 1 & & 0 \\ & & \ddots & \\ 0 & & & 1 \end{bmatrix} \quad (2.13)$$

and we see, that \tilde{e} has a multivariate normal distribution with mean vector

$$E(\tilde{e}) = \frac{b^{-1}}{n} \begin{bmatrix} M \\ 0 \end{bmatrix} = \frac{M}{n} \begin{bmatrix} 1 \\ \phi \\ \vdots \\ \phi^{n-1} \end{bmatrix} \quad (2.14)$$

and variance-covariance matrix

$$\begin{aligned} \hat{I} &= \sigma^2 b^{-1} I b^{-1} = \sigma^2 (b' b)^{-1} = \sigma^2 b^{-1} \\ &= \frac{\sigma^2}{1-\phi^2} \begin{bmatrix} (1-\phi^2) & \phi(1-\phi^2) & \phi^2(1-\phi^2) & \dots & \phi^{n-1}(1-\phi^2) \\ \phi(1-\phi^2) & (1-\phi^4) & \phi(1-\phi^4) & \dots & \phi^{n-2}(1-\phi^4) \\ \phi^2(1-\phi^2) & \phi(1-\phi^4) & (1-\phi^6) & \dots & \phi^{n-3}(1-\phi^6) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \phi^{n-1}(1-\phi^2) & \phi^{n-2}(1-\phi^4) & \dots & \dots & 1-2\phi \end{bmatrix} \end{aligned} \quad (2.15)$$

We see, that for $\phi > 1$ the unconditional expected value of ϵ_1 is growing exponentially (in i), and its variance is also steadily increasing.

In short, the class of models we are concerned with can be written

$$y = X\theta + \epsilon \quad (2.16)$$

with

$$\epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} = \begin{bmatrix} M + \epsilon_1 \\ \phi\epsilon_1 + \epsilon_2 \\ \vdots \\ \phi\epsilon_{n-1} + \epsilon_n \end{bmatrix} \quad (2.17)$$

where ϵ is spherically normally distributed:

$$p(\epsilon) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2}\sigma^2 \epsilon' \epsilon\right). \quad (2.18)$$

If stationarity and reversibility can be assumed a priori, then one degree of freedom may be saved by adopting (2.6) instead of (2.11), i.e. by eliminating the parameter M . In terms of (2.17) this elimination is done by formally replacing M by $((1-\phi)^2)^{-1/2} \epsilon_1$. The resulting subclass of stationary models shall in the following be identified by superscript "o".

2.2 Likelihood functions

In this section we derive the likelihood function of the not necessarily stationary model (in the following referred to as the general model), as well as give that of the a priori stationary model. Also the maximum likelihood estimates (MLE) and the information matrices associated with these two likelihood functions are given.

2.2.1 Derivation of likelihood functions.

From (2.18) and (2.12) we find

$$p(\epsilon) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2}\sigma^2(\epsilon' b b' \epsilon + M(M-2\epsilon_1))\right) \quad (2.19)$$

using, that the Jacobian of the transformation (2.12) is unity:

$$\left| \frac{\partial \tilde{\epsilon}}{\partial \epsilon} \right| = \left| \frac{\partial \begin{bmatrix} b \\ \tilde{\epsilon} \\ 0 \end{bmatrix}}{\partial \epsilon} \right| = \frac{\partial \epsilon' b'}{\partial \epsilon} = |b'| = 1 \quad (2.20)$$

Introducing the matrix

$$C = \begin{bmatrix} -\phi & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \quad (n-1) \times n \quad (2.21)$$

we may write

$$\tilde{\epsilon} = \begin{bmatrix} 0 \\ \tilde{\epsilon} \\ 0 \end{bmatrix} + \begin{bmatrix} e_1' M \\ \tilde{\epsilon} \\ 0 \end{bmatrix} \quad (2.22)$$

or

$$\tilde{\epsilon}' \tilde{\epsilon} = e_1' C \tilde{\epsilon} + (e_1' M)^2 \quad (2.23)$$

where

$$C = \begin{bmatrix} c' & c \\ \tilde{\epsilon}' & \tilde{\epsilon} \end{bmatrix} \quad (2.24)$$

so that (2.19) may be expressed as

$$p(\epsilon) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2}\sigma^2(\tilde{\epsilon}' C \tilde{\epsilon} + (e_1' M)^2)\right). \quad (2.25)$$

Noting that the transformation $\tilde{e} = \tilde{y} - \tilde{X}\tilde{\theta}$ has unit Jacobian we hence have, that the probability density function of \tilde{y} is

$$p(\tilde{y} | \tilde{M}, \tilde{\theta}, \tilde{\sigma}) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2}\sigma^{-2}((\tilde{y}-\tilde{X}\tilde{\theta})'C(\tilde{y}-\tilde{X}\tilde{\theta}) + (\tilde{y}_1 - \tilde{x}_1'\tilde{\theta} - \tilde{M})^2)\right) \quad (2.26)$$

When the observations are considered known, (2.26) is of course also the likelihood function of the general model, and it may be expressed in a particularly simple way in terms of the transformed variables defined by

$$\tilde{z} = \begin{bmatrix} z_2 \\ z_3 \\ \vdots \\ z_n \end{bmatrix} = \tilde{C}\tilde{y} : \tilde{\Xi} = \tilde{C}\tilde{X} \cdot \quad (2.27)$$

(n-1) × 1 (n-1) × p

From

$$(\tilde{y} - \tilde{X}\tilde{\theta})'C(\tilde{y} - \tilde{X}\tilde{\theta}) = (\tilde{z} - \tilde{\Xi}\tilde{\theta})'(\tilde{z} - \tilde{\Xi}\tilde{\theta}) \quad (2.28)$$

it is seen that the likelihood (relative to the original observations \tilde{y}) is

$$i(\tilde{M}, \tilde{\theta}, \tilde{\sigma} | \tilde{y}) = \sigma^{-n} \exp\left(-\frac{1}{2}\sigma^{-2}((\tilde{z} - \tilde{\Xi}\tilde{\theta})'(\tilde{z} - \tilde{\Xi}\tilde{\theta}) + (\tilde{e}_1 - \tilde{M})^2)\right). \quad (2.29)$$

Introducing the further notation

$$SS(\tilde{\theta}, \phi) = (\tilde{z} - \tilde{\Xi}\tilde{\theta})'(\tilde{z} - \tilde{\Xi}\tilde{\theta}) \quad (2.30)$$

the log likelihood of the general model may conveniently be written

$$L(\tilde{M}, \tilde{\theta}, \phi | \tilde{y}) = -n \log \sigma - \frac{1}{2}\sigma^{-2}SS(\tilde{\theta}, \phi) - \frac{1}{2}\sigma^{-2}(\tilde{e}_1 - \tilde{M})^2. \quad (2.31)$$

For the stationary model it is similarly seen (see for example

Sredni [1970]), that the likelihood is

$$i(\tilde{\theta}, \phi | \tilde{y}) = \sigma^{-n} \sqrt{1-\phi^2} \exp\left(-\frac{1}{2}\sigma^{-2}(\tilde{z}^0 - \tilde{\Xi}^0\tilde{\theta})'(\tilde{z} - \tilde{\Xi}^0\tilde{\theta})\right) \quad (2.32)$$

where

$$\tilde{z}^0 = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \tilde{A}\tilde{y} : \tilde{\Xi}^0 = \tilde{A}\tilde{X} \quad (2.33)$$

and using

$$SS^0(\tilde{\theta}, \phi) = (\tilde{z}^0 - \tilde{\Xi}^0\tilde{\theta})'(\tilde{z}^0 - \tilde{\Xi}^0\tilde{\theta}) \quad (2.34)$$

the log likelihood may be expressed as

$$L^0(\tilde{\theta}, \phi, \sigma | \tilde{y}) = -n \log \sigma + \frac{1}{2} \log(1-\phi^2) - \frac{1}{2} \sigma^{-2} SS^0(\tilde{\theta}, \phi) \quad (2.35)$$

2.2.2 Maximum likelihood estimates.

Finding the maximum likelihood estimators $(\hat{\tilde{M}}, \hat{\tilde{\theta}}, \hat{\phi}, \hat{\sigma})$ of the general model, by, as usual, equating the first partial derivatives of L to zero, we find that

$$\frac{\partial L}{\partial \tilde{M}} = 0 \Rightarrow \hat{\tilde{M}} = \tilde{e}_1 = \tilde{y}_1 - \tilde{x}_1'\hat{\tilde{\theta}} \quad (2.36)$$

and subsequently that

$$\frac{\partial L}{\partial \sigma} = 0 \Rightarrow \hat{\sigma}^2 = \frac{SS(\hat{\tilde{\theta}}, \hat{\phi})}{n} \quad (2.37)$$

hence

$$L(\hat{\tilde{M}}, \hat{\tilde{\theta}}, \hat{\phi} | \tilde{y}) = -\frac{n}{2} \log \frac{SS(\hat{\tilde{\theta}}, \hat{\phi})}{n} - \frac{n}{2} \quad (2.38)$$

in other words, the problem of maximizing L with respect to $(\tilde{M}, \tilde{\theta}, \phi, \sigma)$, is reduced to that of minimizing $SS(\tilde{\theta}, \phi)$ with respect to $(\tilde{\theta}, \phi)$.

As far as $\tilde{\theta}$ is concerned this is nothing but an ordinary least squares problem, i.e.

$$\hat{\tilde{\theta}}(\phi) = (\tilde{\Xi}'\tilde{\Xi})^{-1} \tilde{\Xi}'\tilde{z} = (\tilde{X}'\tilde{C}\tilde{X})^{-1} \tilde{X}'\tilde{C}\tilde{y} \quad (2.39)$$

thus

$$SS(\hat{\tilde{\theta}}, \phi) = SS(\hat{\tilde{\theta}}(\phi), \phi) = \tilde{z}'\tilde{\Xi}(\tilde{\Xi}'\tilde{\Xi})^{-1} \tilde{\Xi}'\tilde{z} = \tilde{y}'(\tilde{C}-\tilde{C}\tilde{X}(\tilde{X}'\tilde{C}\tilde{X})^{-1} \tilde{X}'\tilde{C})\tilde{y} \quad (2.40)$$

and

$$SS(\hat{\tilde{\theta}}, \hat{\phi}) = \min_{\tilde{\theta}, \phi} SS(\tilde{\theta}, \phi) = \min_{\phi} SS(\hat{\tilde{\theta}}, \phi) \quad (2.41)$$

The numerical values of $\hat{\phi}$ and $SS(\hat{\tilde{\theta}}, \hat{\phi})$ may be determined using any of the well known search algorithms for locating the minimum of a "well behaved" function.

Noting that

$$SS(\hat{\theta}, \hat{\phi}) = (Y - X\hat{\theta})' C (Y - X\hat{\theta}) = \hat{e}' C \hat{e} \quad (2.42)$$

so that for given $\hat{\theta}$

$$\frac{\partial SS(\hat{\theta}, \hat{\phi})}{\partial \hat{\phi}} = 2 \hat{\phi} \sum_{i=1}^{n-1} \hat{e}_i^2 - 2 \sum_{i=1}^{n-1} \hat{e}_i \hat{e}_{i+1} \quad (2.43)$$

$$= \hat{\phi} = \frac{\sum_{i=1}^{n-1} \hat{e}_i \hat{e}_{i+1}}{\sum_{i=1}^{n-1} \hat{e}_i^2} \quad (2.44)$$

it is seen, that the iterative procedure depicted in Figure 2.1, (see also Cochrane and Orcutt [1949]) produces the MLSE.

For the stationary model it holds that

$$\sigma^2 = \frac{SS^0(\hat{\theta}, \hat{\phi})}{n} \quad (2.45)$$

$$\hat{\theta} = (\hat{\theta}^0, \hat{\theta}^0)' \quad (2.46)$$

and since

$$L^0(\hat{\theta}, \hat{\phi}, \sigma^2 | Y) = -\frac{n}{2} \log SS^0(\hat{\theta}, \hat{\phi}) + \frac{1}{2} \log(1-\hat{\phi}^2) - \frac{n}{2} \quad (2.47)$$

$\hat{\phi}$ may be found from minimizing

$$\left[\frac{SS^0(\hat{\theta}, \hat{\phi})}{(1-\hat{\phi}^2)^{1/n}} \right] = \left[\frac{Y' (A - A X(X'AX)^{-1} X' \hat{\theta}) Y}{(1-\hat{\phi}^2)^{1/n}} \right] \quad (2.48)$$

with respect to $\hat{\phi}$. If the term $\frac{1}{2} \log(1-\hat{\phi}^2)$ in the log likelihood function is neglected an approximate expression similar to (2.44) is found:

$$\hat{\phi} = \frac{\sum_{i=1}^{n-1} \hat{e}_i \hat{e}_{i+1}}{\sum_{i=1}^{n-1} \hat{e}_i^2} \quad (2.49)$$

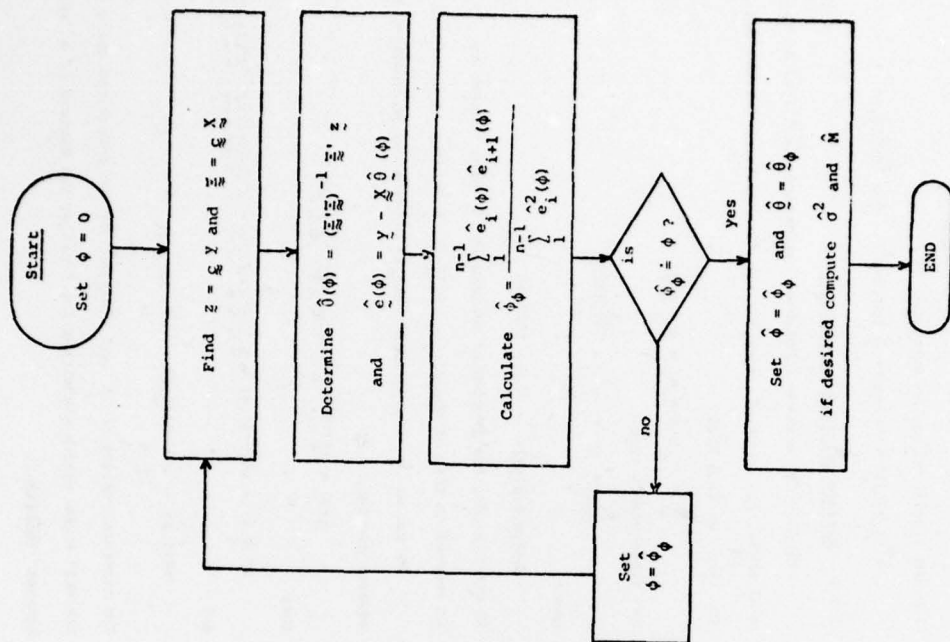


FIGURE 2.1

2.2.3 Fisher's information matrix.

The information matrices J and J^0 corresponding to the log likelihood functions L and L^0 are now given

$$J = \left\{ J_{ij} \right\} = \left[-E \frac{\partial^2 L}{\partial \psi_i \partial \psi_j} \right] \text{ with } \psi_i, \psi_j \in (M, \theta, \phi, \sigma)$$

$$J^0 = \begin{bmatrix} \sigma^{-2} & 0 & 0 & 0 \\ \sigma^{-2} \sum_{i=1}^n (E^0 X_i X_i') & 0 & 0 & 0 \\ 0 & -\sigma^{-2} \sum_{i=1}^n X_i X_i' & \begin{pmatrix} \phi^{2-1} \\ \vdots \\ \phi^{n-1-\phi} n^{-3} \\ \phi^{n-2} \end{pmatrix} J(\phi) \\ 0 & 0 & 0 & 2n\sigma^{-2} \end{bmatrix} \quad (2.50)$$

sym

where

$$J(\phi) = n \left(\frac{1}{1-\phi^2} - \frac{1}{n} \frac{1-\phi^{2n}}{1-\phi^2} + \frac{1}{n} \frac{M^2}{\sigma^2} \frac{1-\phi^{2(n-1)}}{1-\phi^2} \right) \quad (2.51)$$

It is seen, that the information about the starting parameter M does not increase with n ; and that the information about ϕ not only grows with n , but also is inflated when $|\phi|$ is high and when M/σ is large. This is to be expected since under such circumstances the observations will rapidly become very large compared to the error.

About σ it holds, as usual, that the information is proportional to n . As far as θ is concerned we observe, that the information matrix essentially involves the transformed matrix Ξ in the same way as the original matrix X appears in the information matrix corresponding to the usual a priori independent linear model.

The information matrix corresponding to the stationary model is

$$J^0 = \begin{bmatrix} \sigma^{-2} \sum_{i=1}^n \Xi_i \Xi_i' & 0 & 0 \\ 0 & J^0(\phi) & \sigma^{-1} \frac{2\phi}{1-\phi^2} \\ 0 & \sigma^{-1} \frac{2\phi}{1-\phi^2} & 2n\sigma^{-2} \end{bmatrix} \quad (2.52)$$

where

$$J^0(\phi) = n \left(\frac{1}{1-\phi^2} - \frac{1}{n} \frac{1-3\phi^2}{(1-\phi^2)^2} \right). \quad (2.53)$$

As far as θ and σ are concerned this information matrix looks exactly like the one corresponding to the usual a priori independent linear model except the transformed matrix Ξ has taken the place of the original matrix X . The information about ϕ is approximately proportional to n and increases also as $|\phi|$ approaches 1.

2.3 Prior distributions.

In the Bayesian framework the likelihood function does not alone define a statistical model. It must be complemented with a prior distribution for the parameters, and the two components merge into a complete model.

Being part of the model and not a consequence of it, obviously a prior cannot be logically deduced or "proved", it must be chosen. Still the selection of a prior is arbitrary or subjective only to the extent that model specification must always be. The choice of likelihood function is in principle equally arbitrary or subjective.

Of course not all choices of a model to be tentatively entertained are equally defensible, some make more sense than others.

In this section we consider which priors complement the likelihood functions L and L^0 in such a way, that the resulting models seem

intuitively reasonable, and in particular do not lead to clearly unacceptable conclusions. The choice of prior can be interpreted as an expression of prior beliefs concerning the parameters. It may also be regarded as a scaling and shaping of the parameter space in the sense, that the prior chosen for the original parameters ψ may correspond to a locally uniform one in the space of some transformed parameter set $\tilde{\psi}$, so that in this parameter space the normalized likelihood function is equal to the posterior distribution. In so preparing the "stage" for the likelihood, it is kept in mind, that it is not critical to "fine tune" remote regions of the parameter space, where it is barely conceivable that any real data would bring the likelihood function.

Only experience can tell whether the suggested priors are conducive to making the overall model useful in practical data analysis.

2.3.1 Prior for the a priori stationary model.

In general the joint prior distribution of $\{\theta, \phi, \sigma\}$ can be factorized as

$$p^0(\theta, \phi, \sigma) = p^0(\theta, \sigma | \phi) p^0(\phi) \quad (2.54)$$

Focusing attention on the factor $p^0(\theta, \sigma | \phi)$ it is recalled, that for given ϕ the transformed model expression

$$\tilde{z}^0 = \tilde{\Sigma} \tilde{\theta} + \tilde{\epsilon} \quad (2.55)$$

is nothing but a usual independent linear model involving the parameters $\{\theta, \sigma\}$. Following the argument of Box & Tiao [1973] pp. 25-60, a suitable prior for $\{\theta, \sigma\}$ relative to this model form is

$$p_{\phi}^0(\theta, \sigma) \propto \sigma^{-1}. \quad (2.56)$$

The main steps leading to (2.56) are as follows: First an assumption of prior independence between θ and σ is made. Then treating θ and σ separately, Jeffreys' rule is invoked. This rule states, that

for a single parameter, ψ , an approximately noninformative prior distribution results from setting

$$p(\psi) \propto (J(\psi))^{1/2} \quad (2.57)$$

where $J(\psi)$ is Fisher's information measure. The rule has the attractive property, that if $\tilde{\psi}$ is any one to one transformation of ψ , and if $p(\tilde{\psi})$ is also determined from (2.57), then

$$p(\psi)d\psi = p(\tilde{\psi})d\tilde{\psi} \quad (2.58)$$

as is easily seen from the relation

$$J(\tilde{\psi}) = J(\psi) \left(\frac{d\psi}{d\tilde{\psi}} \right)^2. \quad (2.59)$$

Now, if $\tilde{\psi}$ is a transformation of ψ such that $J(\tilde{\psi})$ does not involve $\tilde{\psi}$, then the likelihood is approximately data translated, i.e. different outcomes of $\tilde{\psi}$ have the effect of changing the location of the likelihood function on the $\tilde{\psi}$ axis, while leaving its spread approximately unchanged. In this metric a desire to "let the data speak for themselves" may be expressed by choosing a locally uniform prior for $\tilde{\psi}$. But if $p(\tilde{\psi}) \propto k$ is considered a noninformative prior for $\tilde{\psi}$, then

from (2.58) and (2.59) it follows that a noninformative prior for ψ is

$$p(\psi) = p(\tilde{\psi}) \left| \frac{d\tilde{\psi}}{d\psi} \right| \propto (J(\psi))^{1/2} \quad (2.60)$$

The multiparameter version of Jeffereys' rule has the form

$$p(\tilde{\psi}) \propto |J(\tilde{\psi})|^{1/2} \quad (2.61)$$

If we denote by $\tilde{\psi}$ the parameters in the corresponding metric where $p(\tilde{\psi})$ is locally uniform, then (2.61) ensures that in this metric the "overall spread" as measured by $|J(\tilde{\psi})|^{-1/2}$ remains approximately constant for different outcomes of $\tilde{\psi}$.

As pointed out by Jeffereys [1961], difficulties can occur in applying the multiparameter rule (2.61). In particular when knowledge is available that certain sets of parameters are independent a priori, this must be

included and the rule applied to the independent sets separately.

In the present multiparameter situation we shall proceed as follows.

Assuming prior independence between σ on the one hand and θ on the other we have

$$p^0(\theta, \sigma | \phi) = p^0(\theta | \phi) p^0(\sigma | \phi) \quad (2.62)$$

where $p^0(\theta | \phi)$ is uniform conditionally on ϕ , i.e. must have the form

$$p^0(\theta | \phi) = g(\phi) \quad (2.63)$$

and on the assumption of prior independence between σ and ϕ

$$p^0(\sigma | \phi) = p^0(\sigma) \propto \sigma^{-1} \quad (2.64)$$

It may seem arbitrary to parameterize on σ which is the standard deviation of the innovations (ε_i) , but in fact if one had parameterized on say $\sigma_e = \sigma/\sqrt{1-\phi^2}$, then precisely the same prior would have resulted:

$$p^0(\sigma) = p^0(\sigma_e) \left| \frac{d\sigma_e}{d\sigma} \right| = \left(\frac{\sigma}{\sqrt{1-\phi^2}} \right)^{-1} \frac{1}{\sqrt{1-\phi^2}} = \sigma^{-1} \quad (2.65)$$

Now as far as $p^0(\theta | \phi)$ and $p^0(\sigma)$ are concerned we find

$$p^0(\theta | \phi) p^0(\sigma) = p^0(\theta, \sigma) \propto \left| \varepsilon^0 \varepsilon^0 \right|^{1/2} (1-\phi)^{-1/2} \quad (2.66)$$

using Jeffreys' rule (2.61) jointly for (θ, ϕ) , and where the asymptotic approximation

$$p^0(\phi) \propto (1-\phi^2)^{-1/2} \quad (2.67)$$

has been adopted. It may be noted, that the noninformative prior (2.67) for ϕ , is equivalent to a locally uniform one for $\phi = \arcsin \phi$.

A joint prior distribution for (θ, ϕ, σ) complementing the likelihood function l^0 (2.12) is thus produced by multiplying the two separated independent components (2.64) and (2.67):

$$p^0(\theta, \phi, \sigma) \propto \left| \varepsilon^0 \varepsilon^0 \right|^{1/2} (1-\phi^2)^{-1/2} \sigma^{-1} \quad (2.68)$$

Because multiparameter problems have to be treated with considerable care, it is perhaps valuable to discuss from a somewhat different point of view the prior dependence of θ on ϕ ; i.e. the appropriateness of the factor $p^0(\theta | \phi)$ in the prior.

If a prior is to be selected which does not prejudice the data, it might at first sight seem, that this openmindedness could be expressed by a prior implying, that all values of the vector of expected observations $E(\tilde{y}) = \tilde{\eta}$ are equally acceptable. This notion implies a locally uniform prior for $\tilde{\eta}$. However such a choice is inconsistent with the model structure since it ignores the restrictions (2.2) imposed on $\tilde{\eta}$, as well as it neglects the dependence among the y_i 's.

Multiplying $\tilde{y} = \tilde{\eta} + \tilde{e}$ by the independence inducing matrix \tilde{a} we get

$$\tilde{z}^0 = \tilde{\zeta}^0 + \tilde{\varepsilon} \quad (2.69)$$

where

$$\tilde{\zeta}^0 = \tilde{\varepsilon}^0 \tilde{\theta} \quad (2.70)$$

so that $\tilde{\zeta}^0$ is restricted to the p -dimensional subspace $W(\phi)$ spanned by the column vectors $\tilde{\varepsilon}_1^0, \tilde{\varepsilon}_2^0, \dots, \tilde{\varepsilon}_p^0$ of $\tilde{\varepsilon}^0$. Hence $\tilde{\zeta}^0$ is completely described by its projections onto $W(\phi)$. Let $\tilde{\xi} = [\tilde{\xi}_1, \tilde{\xi}_2, \dots, \tilde{\xi}_p]$ be any orthonormal basis for $W(\phi)$. Now consider the vector

$$\tilde{\theta} = \tilde{\xi}' \tilde{\zeta}^0 \quad (2.71)$$

The j -th element of $\tilde{\theta}$ is the length of the projection of $\tilde{\zeta}^0$ on $\tilde{\xi}_j$. It is now clear, that a locally uniform prior for $\tilde{\theta}$ (conditional on ϕ) is exactly reflecting the desired openmindedness while being consistent with the explicit restrictions dictated by the model structure.

It follows from

$$\tilde{\theta} = \tilde{\xi}' \tilde{\zeta}^0 = \tilde{\xi}' \tilde{\varepsilon}^0 \tilde{\theta} \quad (2.72)$$

that a locally uniform prior for $\tilde{\theta}$ equivalently may be expressed in terms of the following prior for $\tilde{\theta}$:

$$p^0(\tilde{\theta}|\phi) = p^0(\tilde{\theta}) \left| \frac{d\tilde{\theta}}{d\theta} \right| = \left| \xi^0 \cdot \xi \right| \xi^0 \cdot \xi / 2 = \left| \xi^0 \cdot \xi \right|^{1/2} \quad (2.73)$$

which is exactly the same answer as before.

2.3.2 Prior for the not necessarily stationary model.

The joint prior distribution for $(M, \tilde{\theta}, \phi, \sigma)$ may in general be factorized as

$$p(M, \tilde{\theta}, \phi, \sigma) = p(M, \tilde{\theta}|\phi, \sigma) p(\phi) p(\sigma) \quad (2.74)$$

It seems appropriate to assume independence a priori between M and ϕ ; also we shall assume as before that σ is independent a priori of the remaining parameters. On this basis we may write

$$p(M, \tilde{\theta}, \phi, \sigma) = p(M) p(\tilde{\theta}|M, \phi) p(\phi) p(\sigma) \quad (2.75)$$

Following the same line of argument as was put forth in the previous subsection we get

$$p(\sigma) = \sigma^{-1} \quad (2.76)$$

and

$$p(\tilde{\theta}, M|\phi) = p(\tilde{\theta}|M, \phi) p(M) = \left| \xi^0 \cdot \xi \right|^{1/2} \quad (2.77)$$

where clearly the prior for the linear starting parameter M is set locally uniform.

Thus a prior distribution for $(M, \tilde{\theta}, \phi, \sigma)$ complementing the likelihood function l (2.29) has the form

$$p(M, \tilde{\theta}, \phi, \sigma) = p(\phi) \left| \xi^0 \cdot \xi \right|^{1/2} \sigma^{-1} \quad (2.78)$$

where the question of choosing an appropriate prior for ϕ has been left open. This question shall be returned to in Section 2.5.

2.4 Posterior distributions.

At this point Bayes theorem is invoked to establish the posterior density functions of the parameters in the two models, as the normalized product of their respective likelihoods and their complementing priors.

Specifically we have for the general model

$$\begin{aligned} p(M, \tilde{\theta}, \phi, \sigma|\tilde{y}) &= p(M, \tilde{\theta}, \phi, \sigma) l(M, \tilde{\theta}, \phi, \sigma|\tilde{y}) \\ &= p(\phi) \left| \xi^0 \cdot \xi \right|^{1/2} \sigma^{-(n+1)} \exp\left(-\frac{1}{2} \sigma^{-2} (SS(\tilde{\theta}, \phi) + (\epsilon_1 - M)^2)\right) \quad (2.79) \end{aligned}$$

and for the stationary model

$$\begin{aligned} p^0(\tilde{\theta}, \phi, \sigma|\tilde{y}) &= p^0(\tilde{\theta}, \phi, \sigma) l^0(\tilde{\theta}, \phi, \sigma|\tilde{y}) \\ &= \left| \xi^0 \cdot \xi \right|^{1/2} \sigma^{-(n+1)} \exp\left(-\frac{1}{2} \sigma^{-2} SS^0(\tilde{\theta}, \sigma)\right) \quad (2.80) \end{aligned}$$

From a Bayesian point of view, all information which the data \tilde{y} may possess about the parameters $(M, \tilde{\theta}, \phi, \sigma)$ (or alternatively $(\tilde{\theta}, \phi, \sigma)$) is explicitly expressed as the posterior distribution $p(M, \tilde{\theta}, \phi, \sigma|\tilde{y})$ (or $p^0(\tilde{\theta}, \phi, \sigma|\tilde{y})$). And all parametric inference consists of interpreting that distribution.

Usually interest is focused on a subset of parameters in which case the joint distribution is an unnecessary complicated carrier of information. Of course using the Bayesian approach marginal inference is made from considering the relevant marginal densities, which in principle are easily found by integrating out the nuisance parameters.

In the context of the present models the parameters M and σ will typically be of little if any interest by themselves. Accordingly we shall proceed by eliminating M and σ to establish the marginal distributions of $\tilde{\theta}$ and ϕ jointly. In following steps involving further integration, we shall see how inferences are made about $\tilde{\theta}$, or a subset of $\tilde{\theta}$, and about ϕ individually.

2.4.1 Marginal posterior distribution of $\{\hat{\theta}, \hat{\phi}\}$ jointly.

Recognizing, that $p(M, \hat{\theta}, \hat{\phi} | y)$ (2.79) is a Normal function

with respect to M , this parameter may be integrated out to yield

$$p(\hat{\theta}, \hat{\phi} | y) = p(\hat{\theta}) \left| \frac{\partial \hat{\theta}}{\partial \theta} \right| \frac{1}{2} \sigma^{-n} \exp \left(-\frac{1}{2} \sigma^{-2} SS(\hat{\theta}, \hat{\phi}) \right) \int_0^{\infty} \sigma^{-1} \exp \left(-\frac{1}{2} \sigma^{-2} (M - e_1)^2 \right) dM \\ \approx p(\hat{\theta}) \left| \frac{\partial \hat{\theta}}{\partial \theta} \right| \frac{1}{2} \sigma^{-n} \exp \left(-\frac{1}{2} \sigma^{-2} SS(\hat{\theta}, \hat{\phi}) \right) \quad (2.81)$$

This density is of the Gamma form with respect to σ , so integrating over that parameter gives

$$p(\hat{\theta}, \hat{\phi} | y) = p(\hat{\theta}) \left| \frac{\partial \hat{\theta}}{\partial \theta} \right| \frac{1}{2} \int_0^{\infty} \sigma^{-n} \exp \left(-\frac{1}{2} \sigma^{-2} SS(\hat{\theta}, \hat{\phi}) \right) d\sigma \\ \approx p(\hat{\theta}) \left| \frac{\partial \hat{\theta}}{\partial \theta} \right| \frac{1}{2} (SS(\hat{\theta}, \hat{\phi}))^{-\frac{n-1}{2}} \quad (2.82)$$

Turning to the stationary model we find similarly:

$$p^0(\hat{\theta}, \hat{\phi} | y) = \left| \frac{\partial \hat{\theta}}{\partial \theta} \right| \frac{1}{2} (SS^0(\hat{\theta}, \hat{\phi}))^{-n/2} \quad (2.83)$$

2.4.2 Marginal posterior distribution of $\hat{\phi}$.

Integrating over $\hat{\theta}$ we find

$$p(\hat{\phi} | y) = \int_{R^p} p(\hat{\theta}) \left| \frac{\partial \hat{\theta}}{\partial \theta} \right| \frac{1}{2} (SS(\hat{\theta}, \hat{\phi}))^{-\frac{n-1}{2}} d\hat{\theta} \\ = p(\hat{\phi}) \left| \frac{\partial \hat{\phi}}{\partial \phi} \right| \frac{1}{2} \int_{R^p} (SS(\hat{\theta}, \hat{\phi}) + (\hat{\theta} - \hat{\theta})^T \hat{\Sigma}^{-1} (\hat{\theta} - \hat{\theta}))^{-\frac{n-1}{2}} d\hat{\theta} \quad (2.84)$$

Excluding the point(s) $\hat{\phi}^*$ (if any) which makes $\hat{\Sigma}^{-1} \hat{\Sigma}$ singular, this t-form expression integrates to

$$p(\hat{\phi} | y) = p(\hat{\phi}) (SS(\hat{\phi}, \hat{\phi}))^{-\frac{n-p-1}{2}} \quad (2.85)$$

with $-\infty < \hat{\phi} < \infty$.

This exclusion is a matter of formality only since the probability

of $\hat{\phi}$ being equal to the distinct point(s) $\hat{\phi}^*$ is zero (see Section 3.2 of Chapter 3 and Section 4.2 of Chapter 4).

Similarly for the stationary model we find

$$p^0(\hat{\phi} | y) = (SS^0(\hat{\phi}, \hat{\phi}))^{-\frac{n-p}{2}} \quad (2.86)$$

with $-1 < \hat{\phi} < 1$

Attention is at this point directed to Appendix B, where it is studied what the consequences would have been for $p(\hat{\phi} | y)$ had prior independence been assumed between $\hat{\theta}$ and $\hat{\phi}$.

2.4.3 Marginal posterior distribution of $\hat{\theta}$.

The marginal posterior distribution of $\hat{\theta}$ relative to the general model is found as

$$p(\hat{\theta} | y) = \int_{R^p} p(\hat{\theta}, \hat{\phi} | y) d\hat{\phi} \quad (2.87)$$

but here it does not seem possible to perform the integration analytically, so an explicit expression of $p(\hat{\theta} | y)$ in terms of well known functions of y is unavailable.

Utilizing the factorization

$$p(\hat{\theta}, \hat{\phi} | y) = p(\hat{\theta} | \hat{\phi}, y) p(\hat{\phi} | y) \quad (2.88)$$

the integration (2.87) may be carried out conveniently by numerical methods on an electronic computer.

The second factor in (2.88) shall be looked into in Section 2.5 and the first factor is the posterior distribution of $\hat{\theta}$ given $\hat{\phi}$, which from (2.84) is seen to be a multivariate t-distribution, specifically

$$p(\hat{\theta} | \hat{\phi}, y) \sim t_{\hat{\theta}}(\hat{\theta}, s(\hat{\theta}, \hat{\phi})^{-1}, v)$$

with

$$\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2)^T, \quad (2.89)$$

vs $\hat{\phi} = SS(\hat{\phi}, \hat{\phi})$

and

$$v = n - p - 1.$$

Hence substituting (2.88) into (2.87) we see, that $p(\theta|\tilde{y})$ may be regarded as a weighted sum of t-distributions, where $p(\phi|\tilde{y})$ acts as a weight function. In this interpretation the special case of independence is equivalent to setting $p(\phi|\tilde{y})$ equal to an impulse function,

$$\delta(\phi - \phi_0).$$

Of course a numerical integration consists of summing the integrand in small increments of ϕ , multiplied by the increment $\Delta\phi$, i.e. setting

$$p(\theta|\tilde{y}) \approx \Delta\phi \int_{\Delta\phi} p(\theta|\phi, \tilde{y}) p(\phi|\tilde{y}). \quad (2.90)$$

Since $p(\theta|\phi, \tilde{y})$ is known exactly including the normalizing constant, and interpreting $p(\phi|\tilde{y})$ as a weight function, we need not find the normalizing constant of that distribution very precisely in order to assure that $p(\theta|\tilde{y})$ integrates to one. The only requirement is, that the weights, being proportional to $p(\phi|\tilde{y})$, add up to one in the actual summation (2.90). In practice $p(\theta|\tilde{y})$ may often be approximated extremely well by a weighted sum of say 5 t-distributions, corresponding to 5 equidistant ϕ -values covering the interval where the body of $p(\phi|\tilde{y})$ is located.

Frequently it is desired to make inference about a subset of $\tilde{\theta}$. Suppose for example, that the r-dimensional subvector $\tilde{\theta}_1$ is of separate interest, where without loss of generality we may assume

$$\tilde{\theta} = \begin{pmatrix} \tilde{\theta}_1 \\ \tilde{\theta}_2 \end{pmatrix} \quad (2.91)$$

then we find that

$$\begin{aligned} p(\theta_1|\tilde{y}) &= \int_{R^{p-r}} p(\theta|\tilde{y}) d\tilde{\theta}_2 \\ &= \int_R p(\phi|\tilde{y}) \int_{R^{p-r}} p(\theta|\phi, \tilde{y}) d\tilde{\theta}_2 d\phi \\ &= \int_R p(\theta_1|\phi, \tilde{y}) p(\phi|\tilde{y}) d\phi \end{aligned} \quad (2.92)$$

where as before $p(\theta_1|\phi, \tilde{y})$ is a t distribution, specifically

$$p(\theta_1|\phi, \tilde{y}) \sim t_{\tilde{\theta}_1, s^2, D_{11}, v} \quad (2.93)$$

with $\tilde{\theta}_1$, s^2 and v as before, equation (2.89), and where D_{11} comes from

$$(\tilde{\Sigma}^{-1})_{11}^{-1} = \begin{matrix} r & & r \\ & \begin{bmatrix} D_{11} & D_{12} \\ D_{12} & D_{22} \end{bmatrix} & \\ p-r & & p-r \end{matrix} \quad (2.94)$$

So again $p(\theta_1|\tilde{y})$ is a weighted sum of t-distributions conditional on ϕ , the weights being proportional to the relative credibility of the ϕ -values as expressed by the marginal posterior distribution of ϕ .

Exactly parallel results hold obviously for the stationary model. The corresponding formulae relative to this model are generated by substituting n for $n-1$ and inserting superscript "0" wherever applicable.

2.5 Marginal prior and posterior for ϕ in the not necessarily

stationary model.

Turning to the question of choosing a prior for ϕ , two possibilities are considered. One is an approximately non-informative prior $p_n(\phi)$ of the Jeffreys type, the other is a locally uniform one

$$p_u(\phi) \propto k. \quad (2.95)$$

How informative this prior is may be judged by comparison to $p_n(\phi)$, and we shall argue below, that $p_u(\phi)$ is not so unreasonable a choice as it

might first seem. If we adopted this prior (2.95), we get from (2.85):

$$p_u(\phi|y) \propto (SS(\hat{\theta}, \phi))^{-\frac{n-p-1}{2}} \quad (2.96)$$

which incidentally looks much like $p^0(\phi|y)$ (2.86) of the stationary model.

It turns out, that $p_u(\phi|y)$ can be approximated closely by a t -distribution with $n-p-2$ degrees of freedom; which is equivalent to saying that in the vicinity of $\phi = \hat{\phi}$, $SS(\hat{\theta}, \phi)$ is very nearly a quadratic function.

According to Jeffreys' rule (2.60), $p_n(\phi)$ should locally be proportional to $J(\phi)^{1/2}$, where (from the shorthand formula (2.51)):

$$J(\phi) = (n-1) + (n-2)\phi^2 + (n-3)\phi^4 + \dots + 2\phi^{2(n-3)} + \phi^{2(n-2)} + \frac{M^2}{\sigma^2} (1 + \phi^2 + \dots + \phi^{2(n-2)}) \quad (2.97)$$

with $-\infty < \phi < \infty$.

If $|\phi|$ is small ($|\phi| < 1$), then $J(\phi)$ is dominated by the low order terms. In this case we see, that unless the quantity $M^2/(n\sigma^2)$ is approaching unity, $J(\phi)$ may be approximated well by

$$J(\phi) \approx (n-1) + (n-2)\phi^2 + (n-3)\phi^4 + \dots + \phi^{2(n-2)} \quad (2.98)$$

In practice M^2/σ^2 should not often be very large for a truly stationary process. And if that contingency is feared in an experimental situation, any reasonable investigator would intuitively counter by making a comparatively larger number of observations ("allowing the system to settle").

So on this assumption a locally approximately noninformative prior looks like

$$p_n(\phi) \propto ((n-1) + (n-2)\phi^2 + \dots + \phi^{2(n-2)})^{1/2} \quad (2.99)$$

If on the other hand the process is explosive, then the high order terms dominate $J(\phi)$, (2.97). Of course if M^2/σ^2 is small (i.e. the explosion has been monitored from its beginning) then (2.98) is still valid. But even if M^2/σ^2 is large, $p_n(\phi)$ of (2.99) would seem a reasonable approximation, since now a noninformative prior should look like $(1 + \phi^2 + \dots + \phi^{2(n-2)})^{1/2}$ or $(\phi^{2(n-2)})^{1/2}$, which is not much different from (2.99) for large ϕ .

So in any event we shall adopt (2.99) as a locally approximately noninformative prior for ϕ . Figure 2.2 shows $p_n(\phi)$ plotted for $n = 6, 15, 25$ and 60 ; the four curves have been scaled to go through 1 for $\phi = \pm 1$. (The $q(\phi)$ curves are explained below).

Unfortunately a direct multiplication by this noninformative prior in (2.85) does not yield meaningful results. It turns out (see the example in Section 2.6.1) that if it is attempted to derive $p(\phi|y)$ as

$$p_n(\phi) (SS(\hat{\theta}, \phi))^{-\frac{n-p-1}{2}},$$

then this product blows up for large ϕ values (positive or negative), where the likelihood function is approaching zero. The cause of this calamity may be explained by Figures 2.3 and 2.4:

First it is reminded, that the noninformative prior for ϕ is equivalent to a locally uniform prior for some $\phi = \phi(\cdot)$, where the ϕ -metric is such that the likelihood function is approximately data translated. In a single parameter situation the effect of a noninformative prior may be illustrated by drawing likelihood curves in the metric induced by that prior. Of course this is not possible to do for multi-parameter likelihood functions. What can be done, is drawing marginal posterior distributions for ϕ in the original as well as the data

translating metric; this is done in Figures 2.3 and 2.4 for sample sizes $n = 15$ and 25 respectively. Actually an alternative argument for producing a noninformative prior, may consist of determining $\phi(\phi)$ so that the (expected) spread of the posterior distribution of ϕ does not depend on ϕ . This approach is sketched in Appendix A, and as it turns out the very same prior for ϕ is derived that way.

The construction of Figures 2.3 and 2.4 is explained in detail in Appendix A, let it suffice here to point out, that the five solid density curves, $p_u(\phi|y)$, in 2.3a and 2.4a come from (2.96) i.e. correspond to a uniform prior for ϕ ; and they are the results of "data" generated with the true value of ϕ being equal to $-1.2, -.6, .0, .6$ and 1.2 . Clearly the extreme curves ($\phi = \pm 1.2$) are much more concentrated than the ones in between, signifying that in the original metric a change in $\hat{\phi}$ (the mode of these curves) not only shifts the posterior along the ϕ -axis but also changes its shape substantially. In contrast Figures 2.2b and 2.4b show that in the ϕ -metric induced by (2.99) the posterior of ϕ remains comparatively unchanged in spread for various values of $\phi = \phi(\hat{\phi})$. Because the spread of the posterior distribution (or curvature of the likelihood function) not only depends on $\hat{\phi}$ but more generally on the data, sampling variation in the spread of $p(\phi|y)$ is to be expected. The adequacy of the ϕ -metric is supported by the lack of pattern from plot to plot of $p(\phi|y)$ (see also plots in Appendix A); so apparently the prior $p_u(\phi)$ is doing what it is supposed to do. But looking more closely at Figures 2.3b and 2.4b it is observed, that the posteriors $p(\phi|y)$ in the data translating metric, unlike $p_u(\phi|y)$ in Figures 2.3a and 2.4a are not nearly symmetric, but possess heavy tails in one direction or the other. The principle of noninformative priors

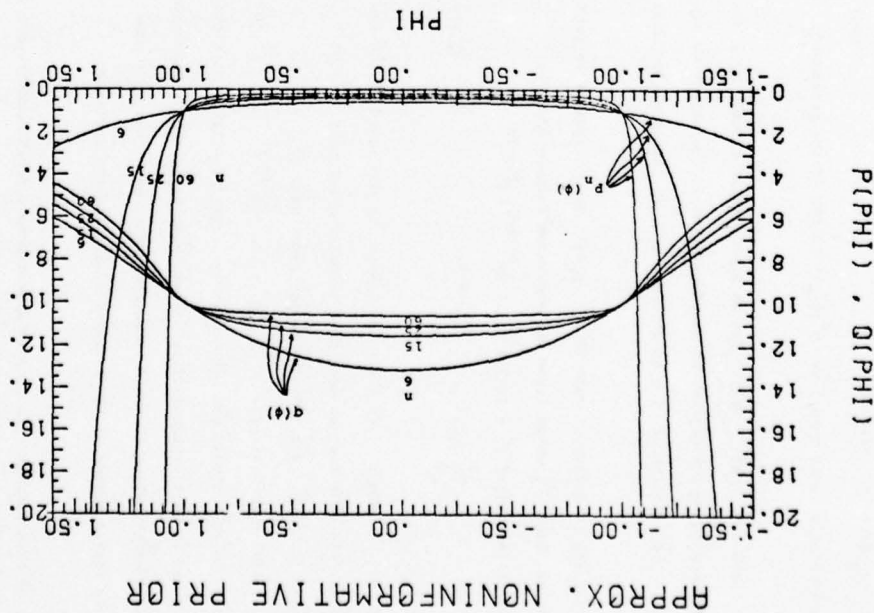


FIGURE 2.2

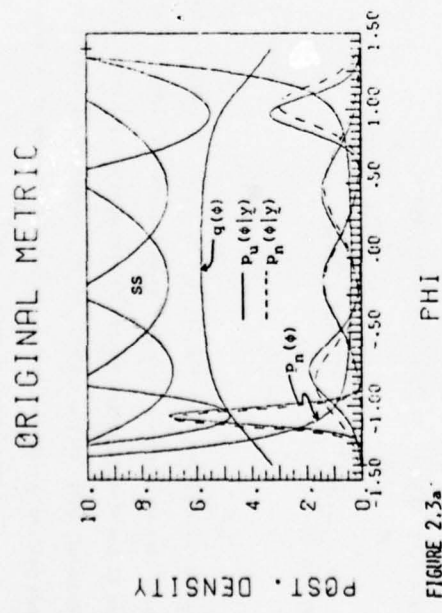


FIGURE 2.3a

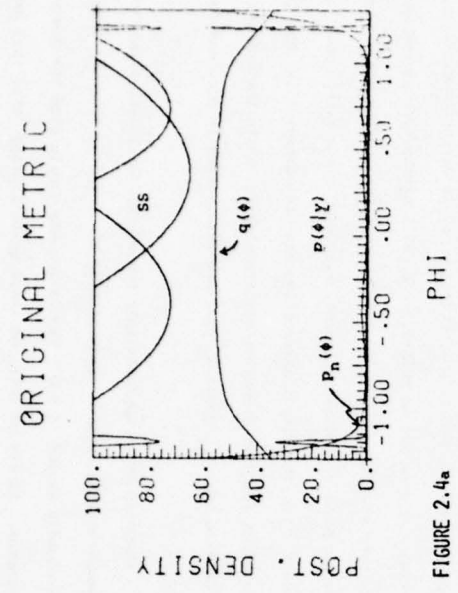


FIGURE 2.4a

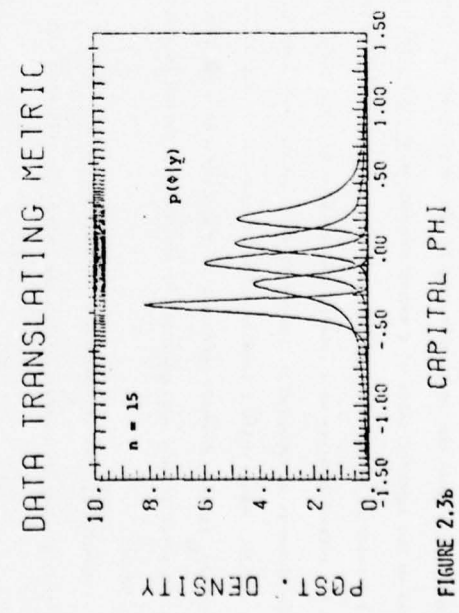


FIGURE 2.3b

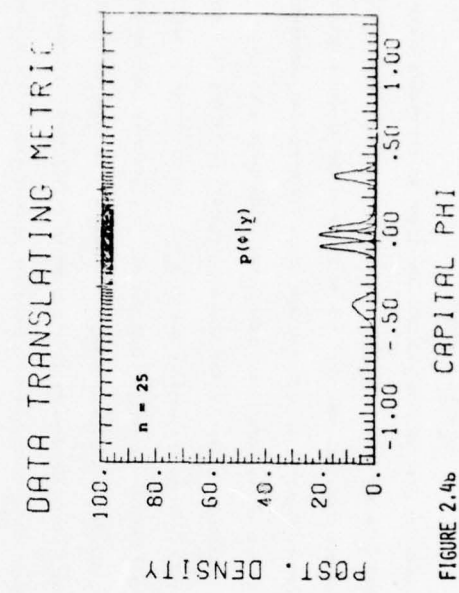


FIGURE 2.4b

does not go beyond the second moment however. Strictly speaking the argument is a local one applicable only in the vicinity of $\hat{\phi}$, where, based on the expected value of a second derivative at this point (used as a measure of spread) the locally approximately noninformative prior has the effect relative to a locally uniform one, of increasing the spread somewhat accompanied by a moderate shift in the mode away from the origin. Hence while a locally uniform prior for ϕ may still be advocated, it is apparently necessary to ensure in some way, that the uniform prior is only applied near $\hat{\phi}$ and does not extend to remote tail areas.

Specifically it is suggested, that the modifying effect of $p_n(\phi)$ on $p_u(\phi|y)$ may be achieved within the t-form the following way. It is noted, that multiplying $SS(\hat{\phi}, \phi)$ by $p_n(\phi)$ is equivalent to multiplying $SS(\hat{\phi}, \phi)$ by $q(\phi)$ in (2.85), where

$$q(\phi) = p_n(\phi) \frac{2}{n-p-1} = \begin{cases} 10 \frac{2}{n-1} \frac{1}{1-\phi^2} \left(1 - \frac{1}{n} \frac{1-\phi^{2n}}{1-\phi^2} \right) - \frac{1}{n-p-1} & \text{when } \phi \neq \pm 1 \\ 10 & \text{when } \phi = \pm 1 \end{cases} \quad (2.100)$$

$q(\phi)$ is plotted in Figure 2.2 for $n = 6, 15, 25, 60$ and $p = 0$; and scaled like the expression (2.100) to go through 10. For $\phi = \pm 1$. Thus if the product $q(\phi) SS(\hat{\phi}, \phi)$ is approximated by a second order polynomial $\bar{Q}(\phi)$ in the vicinity of $\hat{\phi}$, we have the following t-approximation for the marginal posterior distribution of ϕ based locally on the noninformative prior $p_n(\phi)^\dagger$:

[†] The mechanics of this suggested procedures are spelled out in Section 3.4 of Chapter 3.

$$p_n(\phi|y) \propto (Q(\phi))^{-\frac{n-p-1}{2}} \quad (2.101)$$

The posterior distributions (2.101) are drawn as the broken curves in Figures 2.3a and 2.4a. (In the latter figure $p_n(\phi|y)$ and $p_u(\phi|y)$ are almost coincidental even for $\phi = \pm 1.2$ where the dramatic upturn in the prior, Figure 2.2, might perhaps leave the erroneous impression that the effect of the locally noninformative prior grows with n).

Although $p_n(\phi)$ is the advocated (local) prior for ϕ , arguments exist in favor of the uniform prior $p_u(\phi)$. It is simple in application, not only because it is identical for all n , but more importantly because direct multiplication does not create any problems, and therefore an exact posterior density function can be established. Further it may be questioned, whether a noninformative prior really reflects a realistic prior attitude about ϕ . Formally ϕ is defined over the entire interval $(-\infty, \infty)$, but in real life ϕ will hardly ever exceed unity by much. For example a value of $\phi \approx 3$ is rather unthinkable. Such a process would be so violently explosive, that even for a fairly small sample the smallest (i.e. the first) observation and the largest (i.e. the last) observation could be orders of magnitude apart, while the variance of the disturbance ϵ_i added to each new observation is constant (see (2.17)).

Compared to an approximately noninformative prior, a uniform prior is somewhat conservative concentrating the prior probability density more closely around $\phi = 0$. The larger the sample size the greater the difference. It may be felt, at least qualitatively, that this departure reflects a reasonable decreasing willingness to accept larger values of ϕ a priori, as n increases.

Finally it is noted, that $P_n(\phi)$ and $P_u(\phi)$ are much alike for small n , to the point of being identical for the smallest possible sample, $n = 2$. And as n grows large, the likelihood becomes increasingly dominating in the posterior anyway. Hence even if $P_n(\phi)$ is recognized as the superior choice of prior, from the point of view of being unprejudiced relative to the information supplied by the data, the use of $P_u(\phi)$ may still be contemplated, as the practical difference between the two can be expected to be small.

2.6 Examples.

Two examples are presented in this section.

The first example deals with some artificial data, which have been deliberately generated to challenge the method of analysis developed. Admittedly this example is far-fetched, but it does serve to illustrate several points of interest.

The real data of the second example concerned with detecting a possible shift in level of a chemical process variable, are taken from Box and Jenkins [1970].

2.6.1 An artificial example.

These data were generated by the model

$$y_i = \theta_1 x_{i,1} + \theta_2 x_{i,2} + e_i \quad (2.102)$$

with $\theta_1 = 10$, and $\theta_2 = 1$, and where

$$\begin{cases} e_i = M + \epsilon_i \\ e_i = \phi e_{i-1} + \epsilon_i \end{cases} \quad (2.103)$$

with $\phi = .5$, $M = 0$, and the ϵ_i 's are i.i.d. $N(0,1)$.

The values of the independent variable are listed in Table 2.1 along with the synthetic "observations".

TABLE 2.1, The artificial data.

i	y_i	$x_{i,1}$	$x_{i,2}$
1	181.436	1	81.920
2	140.528	1	40.960
3	119.204	1	20.480
4	110.097	1	10.240
5	104.308	1	5.120
6	101.009	1	2.560
7	100.490	1	1.280
8	99.487	1	.640
9	98.130	1	.320
10	98.970	1	.160
11	99.696	1	.080
12	102.054	1	.040
13	100.001	1	.020
14	100.644	1	.010
15	98.955	1	.005

The main purpose of this example is to illustrate a situation where the AR-1 process seems clearly stationary, but where the factors $|\Sigma|^{1/2}$ and $|\Sigma^{(0)}|^{1/2}$ are very different in the priors (2.78) and (2.68) for the general and the stationary model respectively. It is then possible to see how this influences the marginal posterior distribution of ϕ . Also this example demonstrates how the parametric inference about ϕ is affected by the presence of singularity points, of which this data set has two in relation to the general model but only one in relation to the stationary model.

Analyzing this data set on the basis of the general model it is seen, that not only is $\phi^* = 1$ a singularity point because of the mean θ_1 , but also $\phi^* = .5$ is a singularity point because of θ_2 since

$\bar{z}_2 = c \bar{x}_2$ vanishes in this point. This singularity may be explained by recognizing that the exponential decay of $x_{i,2}$, $i = 1, 2, \dots, 15$, can be reproduced in the e_i 's by setting $M = 81.92$ and $\phi = 5$.

How precisely \bar{z} is known given ϕ , depends not only on \bar{v}_2 but also on the $\bar{z}'\bar{z}$ -matrix in that point. Specifically the $100(1-\alpha)\%$ highest posterior density (HPD) region for \bar{z} conditional on ϕ is proportional to $|\bar{z}'\bar{z}|^{-1/2} (vs.)^{1/2}$. The solid $|\bar{z}'\bar{z}|^{-1/2}$ -curve in Figure 2.5a shows how the overall precision with which \bar{z} may be known from the data degrades as ϕ approaches the singularity point ϕ^* .

In this example each singularity is attributed to a single linear parameter, and the two other curves in 2.5a show how the parameters individually become indeterminate as ϕ approaches their respective singularity points. $\sqrt{D_{11}}$ corresponds to $\hat{\theta}_1$ and $\sqrt{D_{22}}$ to $\hat{\theta}_2$, where D_{jj} is defined by Equation (2.94). (All the curves in Figure 2.5a have been scaled to go through 1 at $\phi = 0$).

Figure 2.5b shows posterior distributions for $\phi \cdot P_u(\phi|\bar{y})$ is the exact curve of Equation (2.96). It is noted, that this curve passes smoothly through $\phi = .5$ and $\phi = 1$, and that it can be well approximated by a t -distribution with $n-p-2$ degrees of freedom. The solid curve $P_n(\phi|\bar{y})$ (2.101) has also the t_{n-p-2} form, and is the marginal posterior using locally a noninformative prior as discussed in Section 2.5. Also the result of a direct multiplication of $P_u(\phi|\bar{y})$ by $P_n(\phi)$ is shown, and it is seen, that for $\phi > 1$ where $P_u(\phi|\bar{y})$ is essentially zero, the prior $P_n(\phi)$ begins to overpower the steadily decreasing likelihood and forces the posterior to diverge. This picture is typical for all examples considered, but is only drawn here in Figure 2.5b.

DET FACTOR

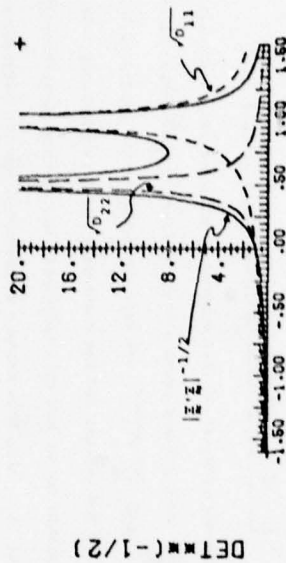


FIGURE 2.5a

MARG. POST. DIST. OF PHI

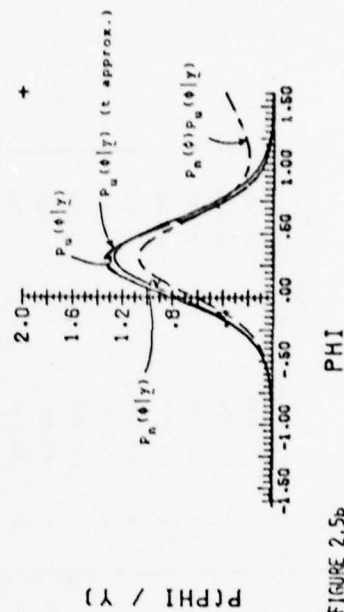


FIGURE 2.5b

Turning to the marginal inference about the linear parameters, the posterior density curves for θ_1 and θ_2 are drawn in Figure 2.6a and Figure 2.6b respectively. The solid curve in each figure is

$$p(\theta_j | \tilde{y}) = \int_R p(\theta_j | \phi, \tilde{y}) p_n(\phi | \tilde{y}) d\phi \quad (2.104)$$

The curve $P_u(\theta_j | \tilde{y})$ is the posterior resulting from applying $P_u(\phi | \tilde{y})$ as a weight function in (2.104) instead of $P_n(\phi | \tilde{y})$. $P_n(\theta_j | \tilde{y})$ and $P_u(\theta_j | \tilde{y})$ are nearly coincident for θ_1 and indistinguishable for θ_2 .

For illustration $P_u(\theta_j | \tilde{y})$ has been approximated, as suggested in Section 2.4.3, by a weighted sum of 5 conditional t-distributions, $p(\theta_j | \phi, \tilde{y})$, specifically those corresponding to $\phi = -.8, -.4, .0, .4, .8$. This "5xt" approximation is very good indeed. It cannot be distinguished from $P_u(\theta_j | \tilde{y})$ for $j = 1$, and deviates only very slightly in the tail for $j = 2$.

Incidentally Figures 2.6a and 2.6b also show the consequences of an independence assumption, $p(\theta_j | \phi=0, \tilde{y})$; as well as the conditional inference on $\phi = \hat{\phi} = .29$. For both θ_1 and θ_2 the failure to recognize the serial correlation has the effect of overestimating the precision with which these parameters are known on the basis of the data, while in this example the location of the various curves are not much different. The approximation $p(\theta_j | \tilde{y}) \approx p(\theta_j | \hat{\phi} = \phi, \tilde{y})$ is good for $j = 1$, but not quite so good for $j = 2$. The probability mass lying outside the depicted interval $.85 < \theta_2 < 1.15$ is 15% for $p(\theta_2 | \tilde{y})$ but only 2% for $p(\theta_2 | \hat{\phi} = \phi, \tilde{y})$.

On the basis of $P_n(\phi | \tilde{y})$ it may be concluded, that there is pretty firm evidence, that the model is stationary, $P(\phi < 1) = 98\%$; so if the further assumption of reversibility is made, then we may reanalyze the data on the basis of the stationary model.

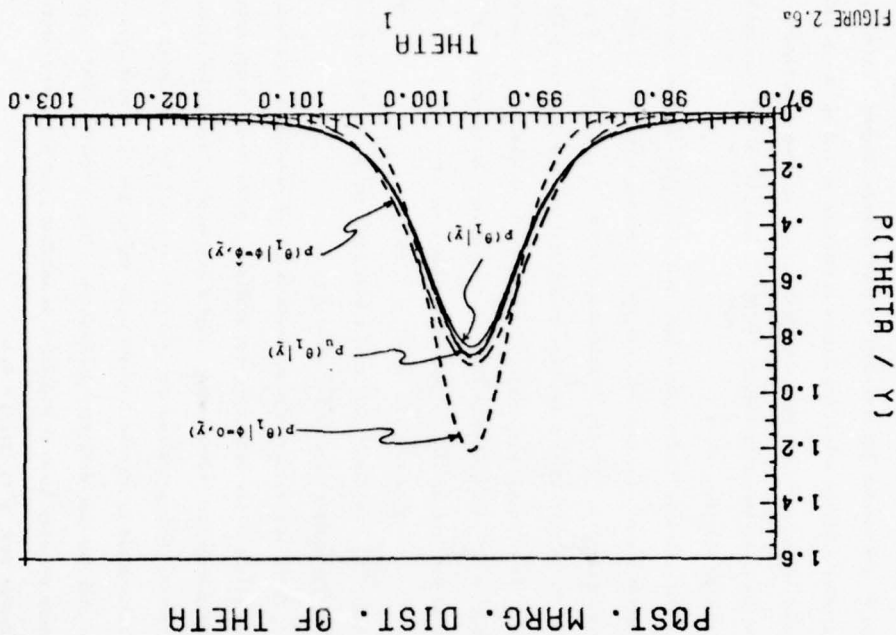


FIGURE 2.6b

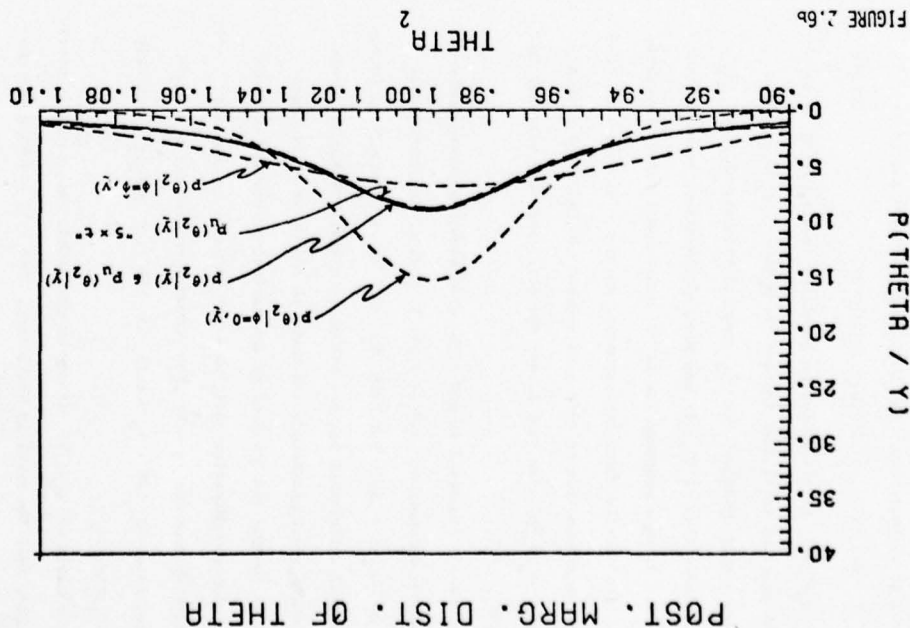


FIGURE 2.6b

The marginal posterior distribution of ϕ , $p^0(\phi|y)$ (2.86), calculated on these premises is drawn in Figure 2.7b. Also $p_u^0(\phi|y)$ and $p_n^0(\phi|y)$ are redrawn for reference in this rescaled figure. Evidently this model revision makes very little difference as far as ϕ is concerned, which is particularly noteworthy since, as seen from Figures 2.5a and 2.7a, the prior factors $|\Xi^0|^{1/2}$ and $|\Xi^{0,0}|^{1/2}$ are quite dissimilar functions of ϕ .

In contrast, Figure 2.8b shows, that $p^0(\theta_2|y)$ is much more concentrated around θ_2 , than was $p(\theta_2|y)$ in Figure 2.6b. The reason for this dramatic increase in information about θ_2 is found in Figure 2.7a. It is seen that in relation to the stationary model, the point $\phi = .5$ is no longer a singularity point, since, unlike Ξ_n^0 , the second column of Ξ^0 does not degenerate to a vector of $(n-1)$ zeros when $\phi = .5$, but (see (2.33)) to the n dimensional vector

$$\Xi_2^0 = (\sqrt{1-.5^2} \ 91.92, 0, \dots, 0)'. \quad (2.105)$$

So now $\sqrt{D_{22}^0}$ is almost flat over a wide range of ϕ values (and substantially smaller than $\sqrt{D_{22}}$ near $\hat{\phi}$).

Under such unusual circumstances it may be beneficial to conduct a reanalysis of the data using the stationary model after having determined through the general model, that a stationarity assumption appears warranted. Ordinarily however $p^0(\theta_j|y)$ and $p(\theta_j|y)$ are quite alike as illustrated by Figures 2.6a and 2.8a where $p(\theta_1|y) = p^0(\theta_1|y)$, and similarly for the conditional posteriors. (Incidentally "5xt" approximations are also drawn in Figures 2.8a and 2.8b, and they are indistinguishable from $p^0(\theta_j|y)$, $j = 1, 2$).

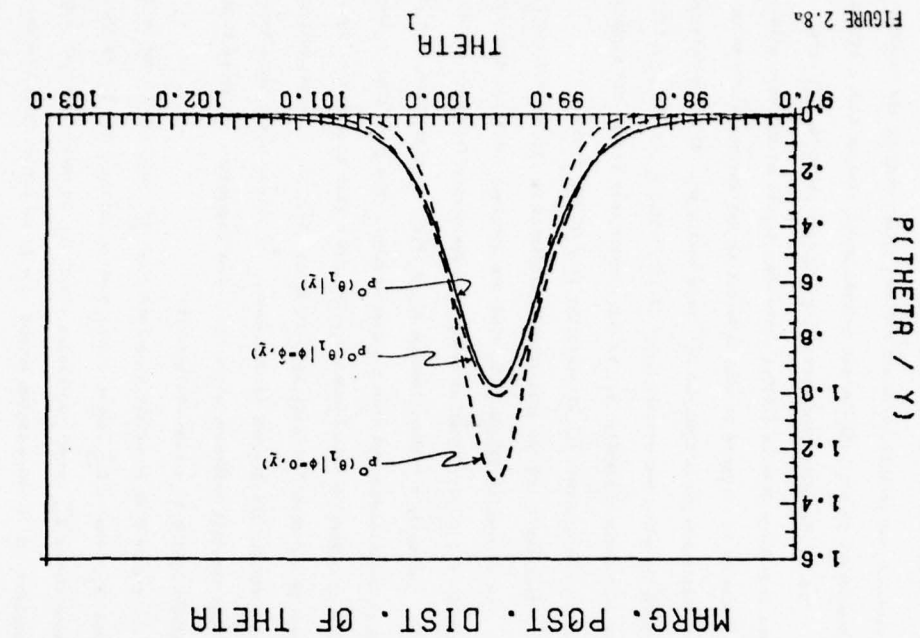


FIGURE 2.7a

MARG. POST. DIST. OF PHI

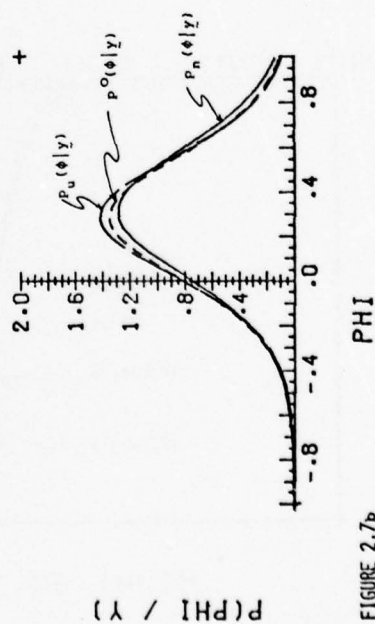


FIGURE 2.7b

MARG. POST. DIST. OF THETA

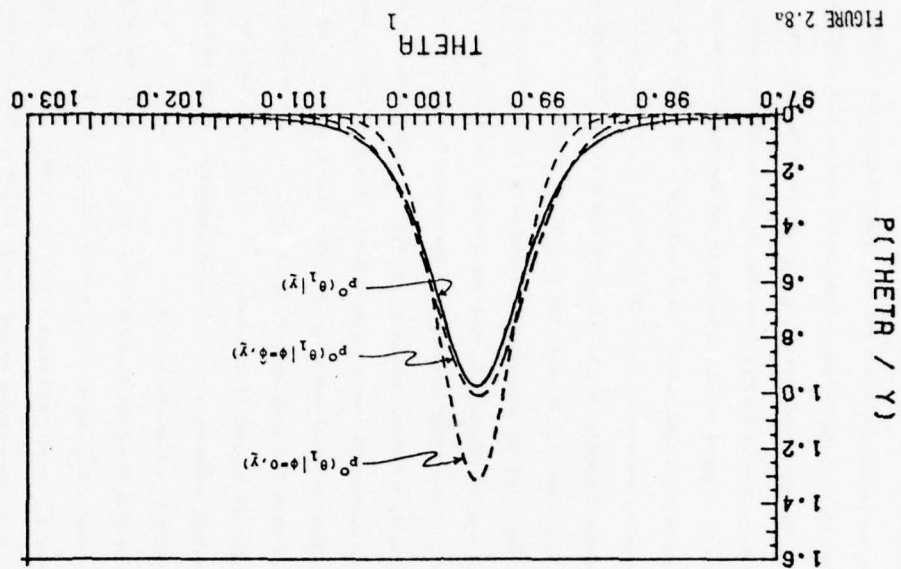


FIGURE 2.8a

2.6.2 Detecting a change in level.

The purpose of this example is to illustrate a situation where ϕ is close to one so that it is not entirely clear whether the process is stationary or not, and how the general model handles such a situation.

The 28 viscosity observations, χ , listed in Table 2.II are taken from Box & Jenkins [1970], where they appear as the last column of their Series D. Suppose it were relevant to ask whether a change in level has taken place from the 22nd observation on. Then the data may be analyzed using the general model (2.16), with $\tilde{\chi}$ as given in Table 2.II; and where evidently θ_1 is the process mean before the suspected change in level, and θ_2 measures the size of the change.

The results of the analysis are displayed in Figures 2.9 and 2.10.

It is seen in Figure 2.9b, that the marginal posterior distribution of ϕ , $p_u(\phi|y)$ goes smoothly from stationary values to explosive ones. (Actually a t-approximation of $p_u(\phi|y)$ is also drawn, but it cannot be distinguished from the exact curve). The posterior $p_n(\phi|y)$ based on a locally noninformative prior has a tail area to the right of one with a probability mass equal to about 15%. So it is not entirely clear whether the process is stationary. Of course the general model allows marginal inference about the linear parameters θ to be made without having to assume stationarity.

Figure 2.9a shows the singularity of $\hat{E}[\tilde{\chi}]$ at $\phi = 1$ due to the mean θ_1 . The $\sqrt{D_{22}}$ curve shows, that in contrast to θ_1 (with its associated $\sqrt{D_{11}}$ curve), inference about θ_2 is apparently not very sensitive to ϕ in a region around $\phi = 1$, as far as spread is concerned. This indication is confirmed by Figures 2.10a and 2.10b showing posterior distributions of θ_1 and θ_2 respectively.

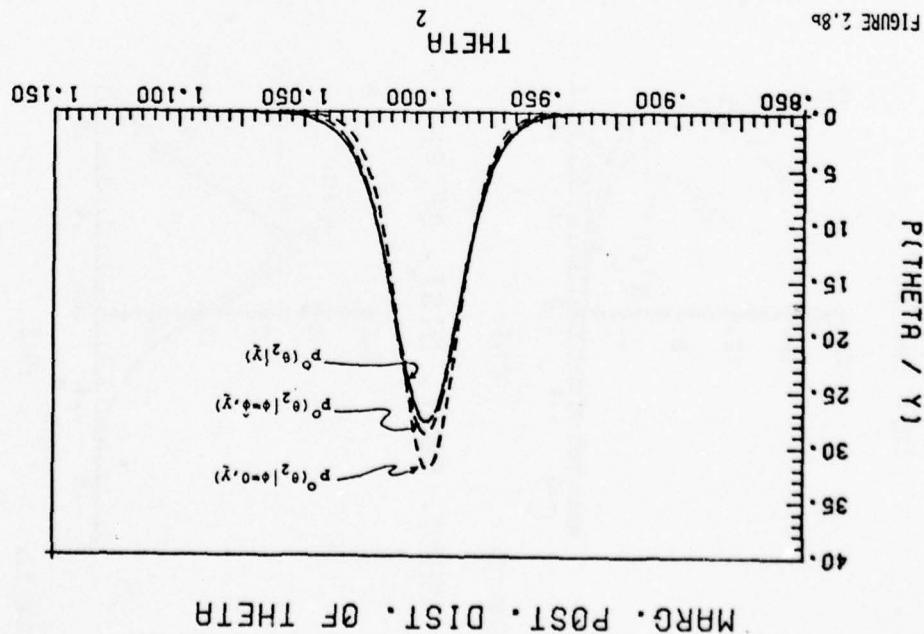


FIGURE 2.8b

DET FACTOR

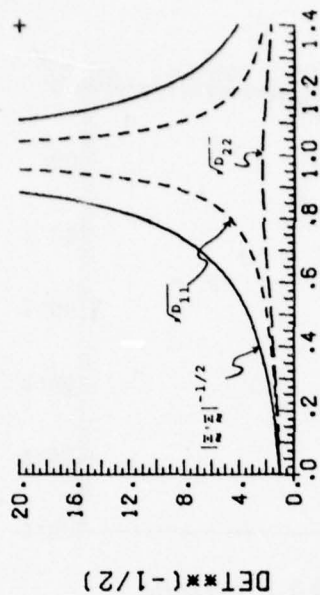


FIGURE 2.9a

MARG. POST. DIST. OF PHI

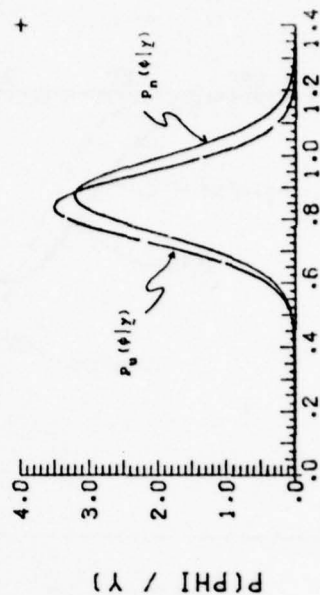
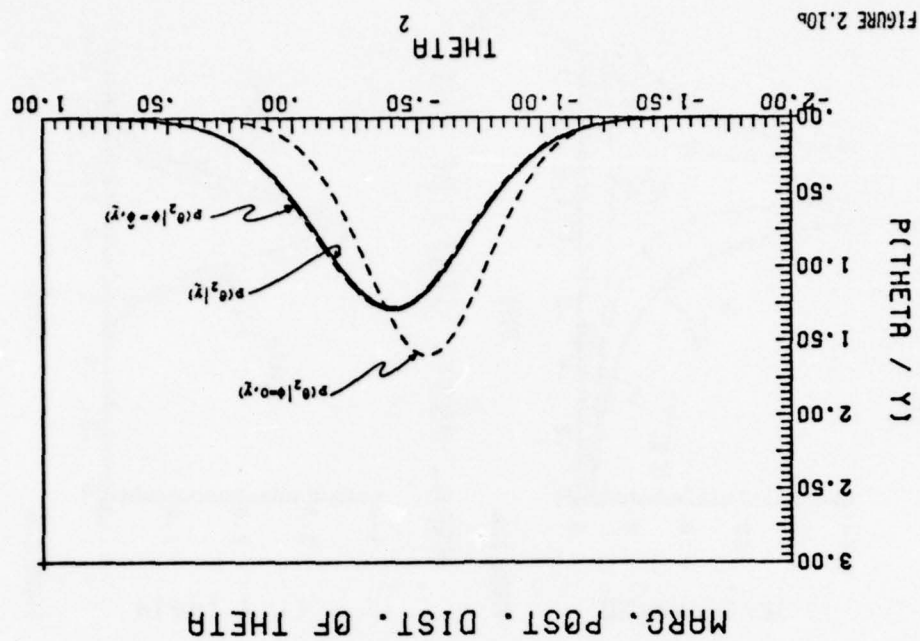


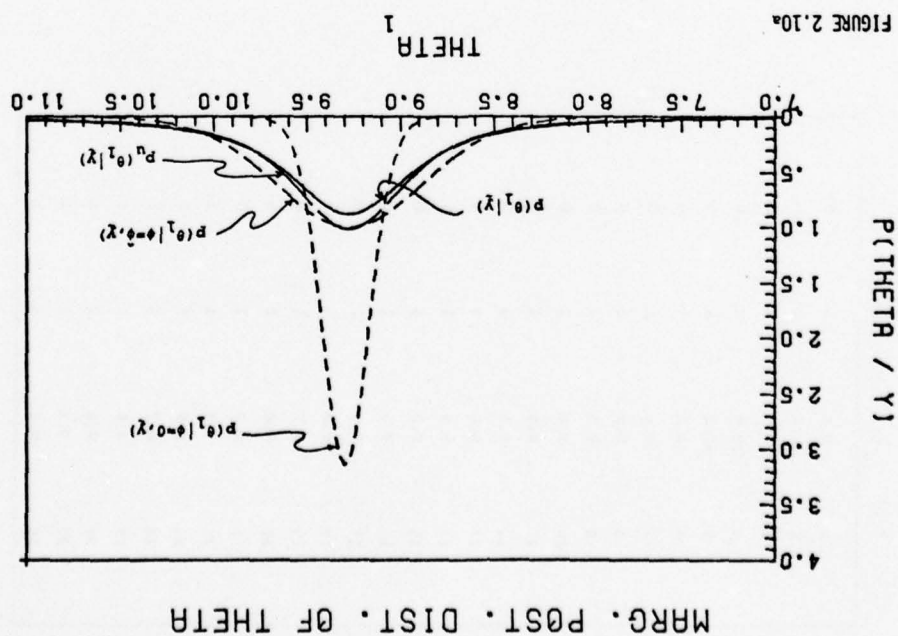
FIGURE 2.9b

TABLE 2.II, The viscosity data.

i	y_i	$x_{i,1}$	$x_{i,2}$
1	9.4	1	0
2	10.0	1	0
3	10.0	1	0
4	10.0	1	0
5	10.2	1	0
6	10.0	1	0
7	10.0	1	0
8	9.6	1	0
9	9.0	1	0
10	9.0	1	0
11	8.6	1	0
12	9.0	1	0
13	9.6	1	0
14	9.6	1	0
15	9.0	1	0
16	9.0	1	0
17	8.9	1	0
18	8.8	1	0
19	8.7	1	0
20	8.6	1	0
21	8.3	1	0
22	7.9	1	1
23	8.5	1	1
24	8.7	1	1
25	8.9	1	1
26	9.1	1	1
27	9.1	1	1
28	9.1	1	1



55



55

MARG. POST. DIST. OF THETA

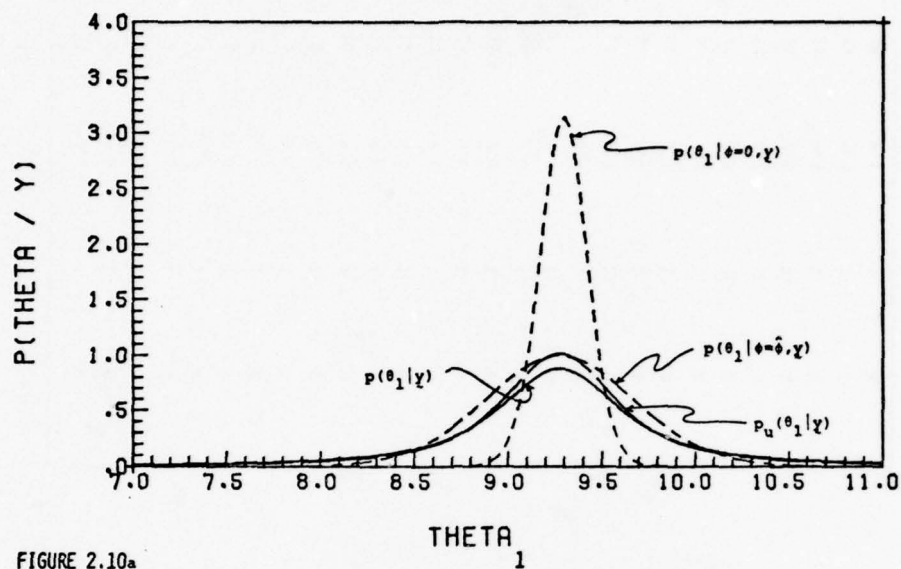


FIGURE 2.10a

55

MARG. POST. DIST. OF THETA

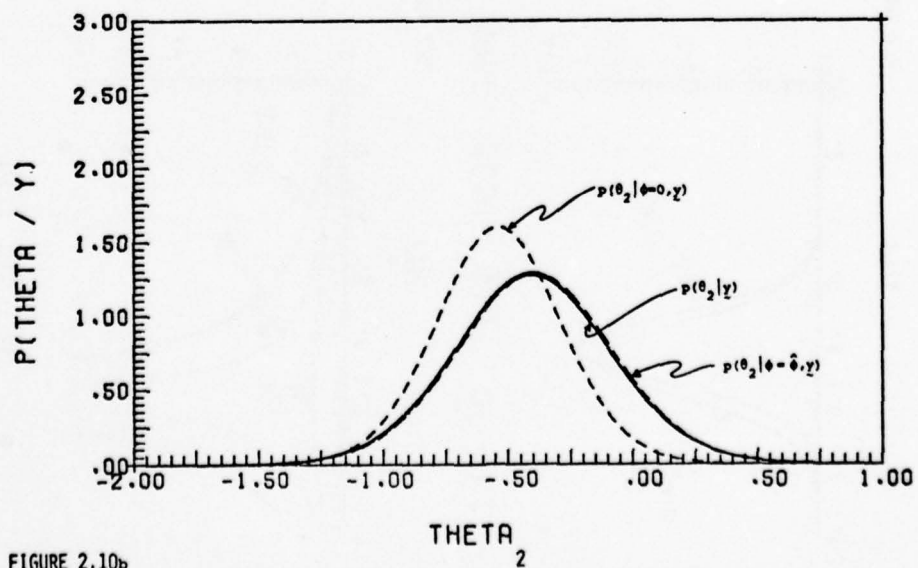


FIGURE 2.10b

56

The marginal posterior $p(\theta_2|y)$ (2.104) in Figure 2.10b has a tail to the right of $\theta_2 = 0$ of 10%, hence the suggested level change is not particularly striking. Incidentally it is noted, that $p_u(\theta_2|y)$, its "5 σ " approximation, and $p(\theta_2|\hat{\phi} = .84, y)$ are virtually identical to $p_n(\theta_2|y)$, while $p(\theta_2|\hat{\phi} = 0, y)$ gives a quite different picture ($P(0, < \theta_2) = 2\%$).

Figure 2.10a shows similarly for the mean θ_1 , that $p_u(\theta_1|y)$ and $p(\theta_1|y)$ are much alike, and a "5 σ " approximation of either curve produces a seemingly perfect match. Also $p(\theta_1|\hat{\phi} = .84, y)$ represents a reasonably good approximation, while an independence assumption, $\hat{\phi} = 0$, yields a very different resulting posterior $p(\theta_1|\hat{\phi} = 0, y)$.

2.7 Conclusion.

The analysis of linear models with independent homoscedastic normal noise occupies a prominent position in statistics. If the assumption of independence cannot be made, it may often hold, that the dependence of the observations can adequately be accounted for by assuming that the noise in the data (or perhaps some transformation of the data, such as e.g. the logarithm or the first difference) follows a first order autoregressive (AR-1) scheme. In this chapter it was developed how such a widened class of linear models may be analyzed from a Bayesian point of view, specifically how inferences can be made about the linear parameters $\hat{\theta}$ and the autoregressive parameter $\hat{\phi}$ jointly, conditionally and marginally.

Two AR-1 schemes were considered. One covers (at a price of one additional starting parameter, M) explosive as well as stationary situations; the other assumes stationarity and reversibility a priori. Their corresponding likelihoods l and l^0 were given, and complemented with

approximately noninformative priors to specify complete models. It was argued, that prior independence between $\hat{\theta}$ and $\hat{\phi}$ is not an appropriate assumption, in particular its adoption may lead to clearly unacceptable conclusions (Appendix B); and it was specifically suggested that $p(\hat{\theta}|\hat{\phi})$ be set proportional to $|\hat{\Sigma}'\hat{\Sigma}|^{1/2} (|\hat{\Sigma}^0\hat{\Sigma}^0|^{1/2})$ for the stationary model) where $\hat{\Sigma}(\hat{\Sigma}^0)$ is the matrix of transformed independent variables, which depends on $\hat{\phi}$.

In situations where the data clearly have stationary noise the marginal posterior distribution for $\hat{\phi}$, $p(\hat{\phi}|y)$, appears very little affected by which of the two AR-1 schemes is employed, even in a case where $|\hat{\Sigma}'\hat{\Sigma}|^{1/2}$ and $|\hat{\Sigma}^0\hat{\Sigma}^0|^{1/2}$ were drastically different. This incidentally supplies further support for the appropriateness of this factor in the prior. For truly stationary situations $\hat{\Sigma}'\hat{\Sigma}$ and $\hat{\Sigma}^0\hat{\Sigma}^0$ will ordinarily be quite alike, so that little is to be gained from employing the a priori stationary model even in such cases. The not necessarily stationary model offers greater flexibility not only to let the data speak to the possibility of explosive noise, but also allows the disturbance associated with the first observation to be untypical. For this more general model, it was discovered, that special care has to be exercised in applying the locally approximately noninformative prior for $\hat{\phi}$. A procedure for overcoming this difficulty was suggested, but as far as making inference about $\hat{\theta}$ it appears from the examples, that it makes virtually no difference if instead a convenient locally uniform prior is utilized. In fact since prior information normally exists that $\hat{\phi}$ cannot be much larger than 1 in practice, an informative prior like the uniform one may make good sense, because it concentrates the prior density more closely around $\hat{\phi} = 0$ relative to a

noninformative prior. In any event the posterior density of ϕ derived from a uniform prior, can always be multiplied by any other prior to give its resulting posterior distribution.

The numerical integration leading to the marginal posterior distribution of θ , $p(\theta|y)$ may be carried with a fine grid for ϕ . It was demonstrated however, that excellent approximations result from computing $p(\theta_j|y)$ as a weighted sum (the weights being proportional to $p(\phi_j|y)$) of five t-distributions, $p(\theta_j|\phi_j, y)$, corresponding to five equally spaced ϕ values covering the interval where the body of $p(\phi|y)$ is located.

Appendix A. On the noninformative prior for ϕ .

In Section 2.3 approximately noninformative priors were produced on the basis of a data translating likelihood argument. As mentioned in Section 2.5 an alternative principle may consist of determining the prior $p(\psi)$ for ψ which corresponds to a locally uniform prior for ψ , where the ψ -metric is such that the posterior density of ψ is approximately data translated.

In the following, this approach is outlined as it applies to the parameter ϕ in the general model. As it turns out, the very same (local) prior results, hence Figures A.1 through A.6 (plus 2.3 and 2.4) showing approximate data translated posteriors in the ϕ -metric equally illustrate the effect of adopting the data translating likelihood principle.

From (2.85) we have (writing $SS(\phi)$ for $SS(\hat{\theta}, \phi)$):

$$p(\phi|y) = (SS(\phi))^{-\frac{v+1}{2}} \quad (A.1)$$

$$\text{where } \phi = \hat{\phi}(\phi) = f(\phi) \quad (A.2)$$

$$\text{and } v = n-p-2. \quad (A.3)$$

In this formulation, we shall determine $\phi = f^{-1}(\phi)$ such that the second derivative of (A.1), is approximately constant, in the point where the first derivative is zero. Actually f^{-1} (or f) need not be determined; it is seen from

$$p(\phi) \propto \frac{d\phi}{d\hat{\phi}} = \frac{df^{-1}(\phi)}{d\phi} \quad (A.4)$$

that the first derivative of f^{-1} (or of f) will suffice.

Because the posterior (A.1) is only specified up to a multiplicative constant, the log-posterior (lp for short) shall be employed, i.e. $p(\phi)$ is to be determined from

$$\frac{\partial^2 L_P}{\partial \phi^2} \bigg|_{\frac{\partial L_P}{\partial \phi} = 0} = \text{constant} \quad (\text{A.5})$$

First

$$\frac{\partial L_P}{\partial \phi} = -\frac{\partial L_P}{\partial \phi} \frac{d\phi}{d\hat{\phi}} = -\frac{v+1}{2} \frac{\partial \ln SS(\hat{\phi})}{\partial \hat{\phi}} f'(\hat{\phi}) \quad (\text{A.6})$$

Equating (A.6) to zero determines $\hat{\phi}$ and $\hat{\phi} = f^{-1}(\hat{\phi})$.

Next

$$\begin{aligned} \frac{\partial^2 L_P}{\partial \phi^2} \bigg|_{\frac{\partial L_P}{\partial \phi} = 0} &= -\frac{v+1}{2} \frac{\partial^2 \ln SS(\hat{\phi})}{\partial \hat{\phi}^2} \left(\frac{d\hat{\phi}}{d\phi} \right)^2 \bigg|_{\hat{\phi}} \\ &= -\frac{v+1}{2} \left(\frac{d\hat{\phi}}{d\phi} \right)^2 \left[\frac{1}{SS(\hat{\phi})} \frac{\partial^2 SS(\hat{\phi})}{\partial \hat{\phi}^2} - \left(\frac{\partial \ln SS(\hat{\phi})}{\partial \hat{\phi}} \right)^2 \right] \bigg|_{\hat{\phi}} \\ &= -\frac{v+1}{2} \left(\frac{d\hat{\phi}}{d\phi} \right)^2 \frac{1}{SS(\hat{\phi})} \frac{\partial^2 SS(\hat{\phi})}{\partial \hat{\phi}^2} \bigg|_{\hat{\phi}} \end{aligned} \quad (\text{A.7})$$

Thus, to satisfy (A.5) it is required that

$$\left(\frac{d\hat{\phi}}{d\phi} \right)^2 = \frac{1}{SS(\hat{\phi})} \frac{\partial^2 SS(\hat{\phi})}{\partial \hat{\phi}^2} \bigg|_{\hat{\phi}} \quad (\text{A.8})$$

Since the right hand side of (A.8) not only depends on the data through $\hat{\phi}$, we shall settle for the lesser requirement, that the expected value of (A.5) is nearly constant. From

$$E \left(\frac{1}{SS(\hat{\phi})} \right) = \frac{1}{v_0^2} \quad (\text{A.9})$$

and

$$\begin{aligned} E \left(\frac{\partial^2 SS(\hat{\phi})}{\partial \hat{\phi}^2} \right) &= E \left[\sum_{i=1}^{n-1} e_i^2 \right] = 2 \sum_{i=1}^{n-1} E \left[\sum_{k=1}^{n-1} e_i^2 \right] \\ &= 2 \sum_{i=1}^{n-1} \left(\sum_{k=1}^{n-1} \phi^{2(i-1)} M^2 + \sum_{k=1}^{n-1} \phi^{2(k-1)} \sigma^2 \right) \end{aligned} \quad (\text{A.10})$$

we find

$$\begin{aligned} \left(\frac{\partial f^{-1}(\hat{\phi})}{\partial \hat{\phi}} \right)^2 &= \left(\frac{d\hat{\phi}}{d\phi} \right)^2 \\ &= \frac{1}{(n-p-2)\sigma^2} \left(n \frac{\sigma^2}{1-\phi} - \frac{\sigma^2 (1-\phi)^{2n}}{(1-\phi)^2} + M^2 \frac{1-\phi^{2(n-1)}}{1-\phi^2} \right) \\ &= \frac{n}{n-p-2} \beta(\hat{\phi}) \end{aligned} \quad (\text{A.11})$$

So as before we are led to the locally noninformative prior

$$p(\phi) = \left\{ \frac{1}{1-\phi} - \frac{1}{n} \frac{1-\phi^{2n}}{(1-\phi)^2} \right\}^{1/2} \quad (\text{A.12})$$

Although f^{-1} (or f) is not found analytically the data translating metric induced by the prior (A.12) can be constructed numerically. This is done in Figures A.1 through A.6 as well as in Figures 2.3 and 2.4. The "data" used in these 8 figures were generated as follows:

For Figures A.1, A.2, 2.3 and 2.4 the generating model is

$$Y_i = e_i \quad i = 1, 2, \dots, n \quad (\text{A.13})$$

with

$$\begin{cases} e_1 = \epsilon_1 \\ e_i = \phi e_{i-1} + \epsilon_i \end{cases} \quad (\text{A.14})$$

where the ϵ_i 's are i.i.d. $N(0,1)$. For the remaining figures the generating models are

$$\begin{aligned} \text{Figure A.3: } Y_i &= \theta_1 + e_i \\ &\text{with } \theta_1 = 1. \end{aligned} \quad (\text{A.15})$$

$$\begin{aligned} \text{Figure A.4: } Y_i &= \theta_1 + \theta_2 i + e_i \\ &\text{with } \theta_1 = 1. \text{ and } \theta_2 = 1. \end{aligned} \quad (\text{A.16})$$

$$\begin{aligned} \text{Figure A.5: } Y_i &= \theta_1 + \theta_2 2^{i-1} + e_i \\ &\text{with } \theta_1 = 1. \text{ and } \theta_2 = .01 \end{aligned} \quad (\text{A.17})$$

(i.e. a model like the one considered in Section (2.6.1)).

$\phi = -1.2, -.6, .0, .6, 1.2$. Looking at the "a"-figures relating to the original metric, the solid (nearly-)parabolas in the upper part of the figure are (scaled) SS-curves. Next the U-shaped prior $p_n(\phi)$ and the corresponding $q(\phi)$ factor are recognized (see Figure 2.2). The solid distribution curves $p_u(\phi|y)$ are those resulting from adopting a locally uniform prior, while the broken curves are posteriors resulting from a locally noninformative prior. The ϕ -metric induced by this prior is depicted in the "b"-figures. The divisions along the upper horizontal axis give a visual impression of how the original scale is stretched and compressed relative to the ϕ -metric. The five posterior distribution curves of ϕ are those given by (A.1).

Figure A.6: $y_i = \theta_1 + \theta_2 x_i + e_i$ (A.18)

with $\theta_1 = 1$ and $\theta_2 = 1$.

and where x_i is a random walk $x_i = \sum_{j=1}^i a_j$,

the a_j 's being i.i.d. $N(0,1)$

(i.e. a model like Model A Equation (3.67) in Section 3.6 of

Chapter 3.)

For all figures the employed sample size n was 15, except for Figures 2.4 and A.2 where $n = 25$. (So differences between Figures 2.3 and A.1 are due to sampling variation only, and the same is true for Figure 2.4 vs. Figure A.2).

A reason for considering the model (A.13) is that the factor

$|\Sigma^{-1}|^{1/2}$ in the prior vanished in this case, which perhaps makes the effect of the noninformative prior for ϕ more transparent. Also the posterior distribution of ϕ , $p_\phi(\phi|y)$, in relation to (A.13) is equal to the posterior of ϕ relating to any model

$$\tilde{y} = X\theta + e \quad (A.19)$$

conditional on $\tilde{\theta}$:

$$p_\phi(\phi|\tilde{y}) = p_n(\phi)(SS(\phi))^{-\frac{n-1}{2}} \quad (A.20)$$

$$= p_n(\phi)(e'Ce)^{-\frac{n-1}{2}} = p(\phi|\tilde{\theta}, \tilde{y})$$

For all the other models we have

$$p(\phi|\tilde{y}) = p_n(\phi)(SS(\tilde{\theta}, \phi))^{-\frac{n-p-1}{2}} \quad (A.21)$$

The various graphs in Figures 2.3, 2.4 (in Chapter 2) and A.1 through A.6 shall now be identified. Each figure depicts posterior distributions corresponding to five "data sets", all five generated from the same e_i 's (different from figure to figure), but with

ORIGINAL METRIC

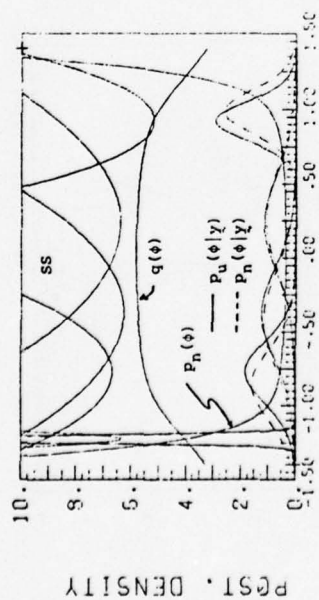


FIGURE A.1a

PHI

DATA TRANSLATING METRIC

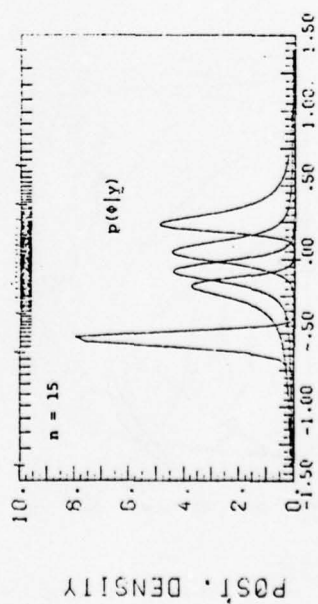


FIGURE A.1b

CAPITAL PHI

ORIGINAL METRIC

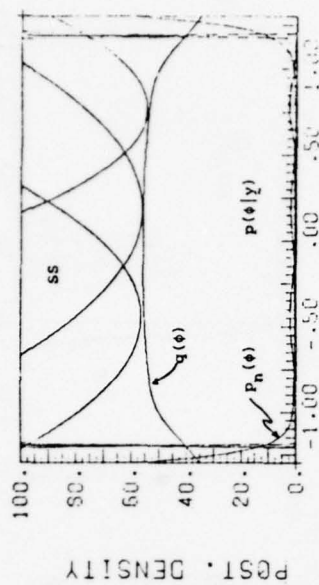


FIGURE A.2a

PHI

DATA TRANSLATING METRIC

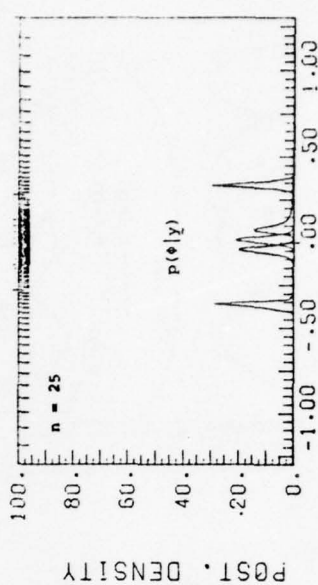


FIGURE A.2b

CAPITAL PHI

ORIGINAL METRIC

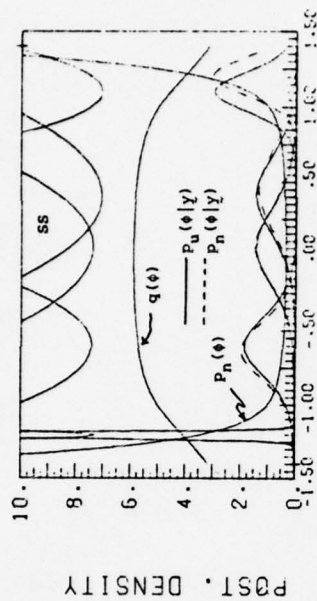


FIGURE A.4a

DATA TRANSLATING METRIC

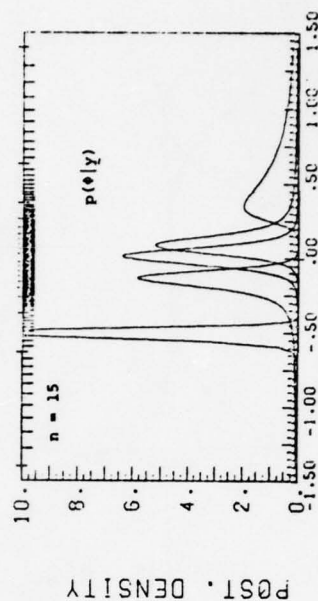


FIGURE A.4a

ORIGINAL METRIC

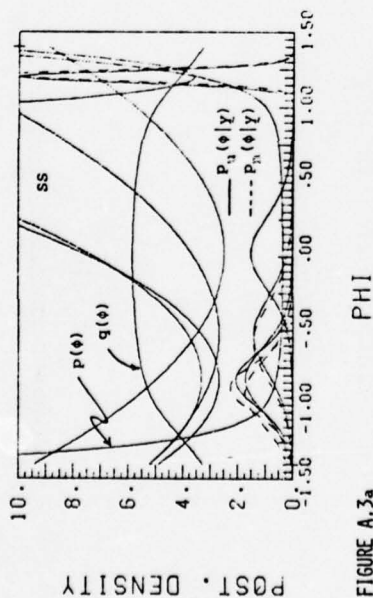


FIGURE A.3a

DATA TRANSLATING METRIC

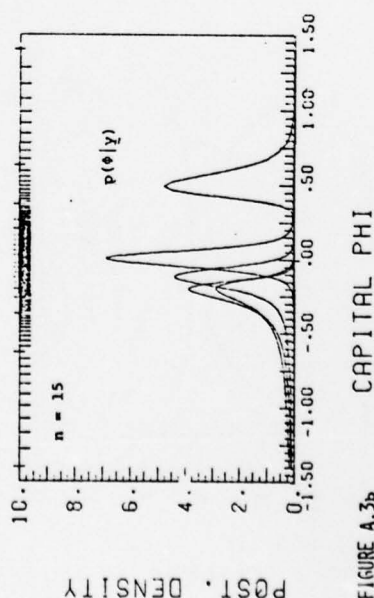


FIGURE A.3b

ORIGINAL METRIC

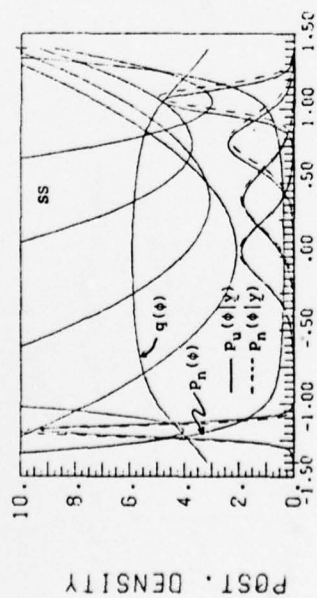


FIGURE A.5a

DATA TRANSLATING METRIC

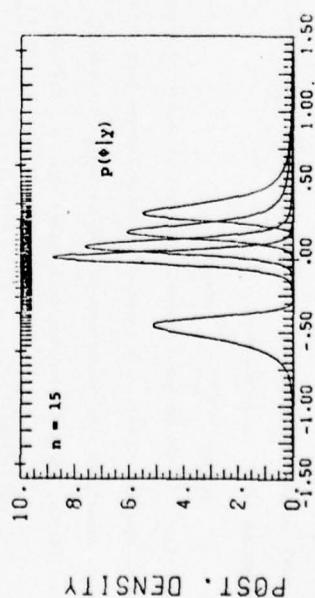


FIGURE A.5b

CAPITAL PHI

ORIGINAL METRIC

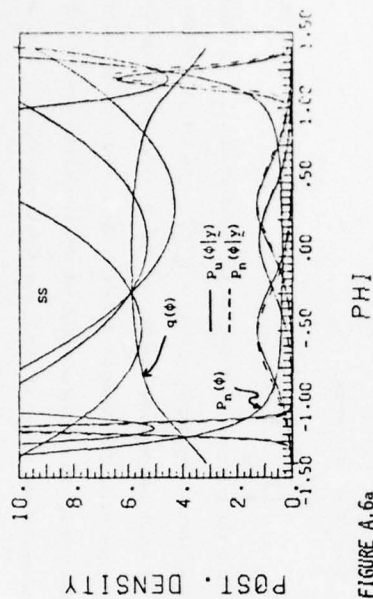


FIGURE A.6a

DATA TRANSLATING METRIC

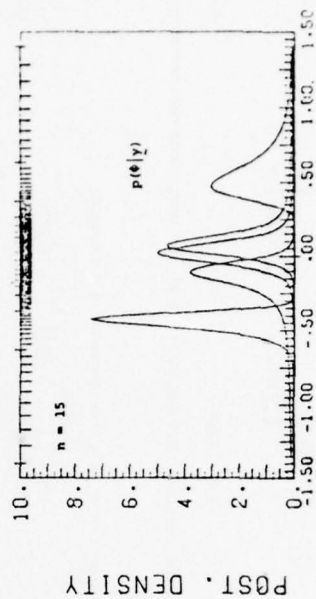


FIGURE A.6b

CAPITAL PHI

Appendix B. Consequences of assuming prior independence between $\underline{\theta}$ and $\underline{\phi}$.

If independence is assumed between $\underline{\theta}$ and $\underline{\phi}$ a priori, then the factor $p(\underline{\theta}|\underline{y})$ becomes a constant in the joint prior for the parameters $(\underline{\theta}, \underline{\phi}, \sigma)$ of the general model. Consequently the marginal posterior distribution of $\underline{\phi}$ takes the form

$$p_1(\underline{\phi}|\underline{y}) = |\underline{\Xi}'\underline{\Xi}|^{-1/2} p(\underline{\phi}|\underline{y}) \quad (\text{B.1})$$

where $p(\underline{\phi}|\underline{y})$ is the marginal posterior derived in Section 2.4 and given in Equation (2.85). The question we shall address here is how $p_1(\underline{\phi}|\underline{y})$ behaves in a neighbourhood of (the excluded) singularity points $\underline{\phi} = \underline{\phi}^*$ if any, where $\underline{\Xi}$ has not full column rank, p . Writing

$$\underline{\Xi} = [\underline{\Xi}'\underline{\Xi}]^{-1/2} \underline{\Xi}'\underline{\Xi} \quad (\text{B.2})$$

it may be assumed without loss of generality (see Section 3.2 of Chapter 3) that $\underline{\Xi}_1 \rightarrow 0$ for $\underline{\phi} \rightarrow \underline{\phi}^*$, while $[\underline{\Xi}'\underline{\Xi}]^{-1/2}$ has full rank p in a neighbourhood of $\underline{\phi}^*$, where $\underline{\Xi}_1$ is the first column of $\underline{\Xi}$ excluding the last row.

From (2.27) we have

$$\underline{\Xi}_1 = \begin{bmatrix} x_{2,1} - \underline{\phi} x_{1,1} \\ \vdots \\ x_{n,1} - \underline{\phi} x_{n-1,1} \end{bmatrix} = \begin{bmatrix} (\underline{\phi}^* - \underline{\phi}) x_{1,1} \\ \vdots \\ (\underline{\phi}^* - \underline{\phi}) x_{n-1,1} \end{bmatrix} = \underline{\delta} \underline{x}_1 \quad (\text{B.3})$$

as $\underline{\phi}^* = x_{2,1}/x_{1,1} = \dots = x_{n,1}/x_{n-1,1}$, and where

$$\underline{\delta} = \underline{\phi}^* - \underline{\phi} \quad (\text{B.4})$$

Now

$$\begin{aligned} \det(\underline{\Xi}'\underline{\Xi}) &= \det([\underline{\Xi}_1' \underline{\Xi}_1 + \underline{\Xi}'(1) \underline{\Xi}(1)]^{-1/2} [\underline{\Xi}_1' \underline{\Xi}_1 + \underline{\Xi}'(1) \underline{\Xi}(1)]) \\ &= \delta^2 \det \begin{bmatrix} \underline{\Xi}_1' \underline{\Xi}_1 & \underline{\Xi}_1' \underline{\Xi}(1) \\ \underline{\Xi}'(1) \underline{\Xi}_1 & \underline{\Xi}'(1) \underline{\Xi}(1) \end{bmatrix} \\ &= \delta^2 d(\underline{\phi}) \end{aligned} \quad (\text{B.5})$$

where the function (determinant) $d(\underline{\phi})$ is continuous and not zero at $\underline{\phi} = \underline{\phi}^*$.

For $0 < \epsilon_1 < \epsilon_2$ we consider the integral

$$\begin{aligned} \lim_{\epsilon_1 \rightarrow 0} \int_{\underline{\phi}^* - \epsilon_2}^{\underline{\phi}^* - \epsilon_1} |\underline{\Xi}'\underline{\Xi}|^{-1/2} p(\underline{\phi}|\underline{y}) d\underline{\phi} \\ = \lim_{\epsilon_1 \rightarrow 0} (d(\underline{\phi})^{-1/2} p(\underline{\phi} = \underline{\phi}_d|\underline{y}) \int_{\underline{\phi}^* - \epsilon_2}^{\underline{\phi}^* - \epsilon_1} \frac{1}{\delta} d\underline{\phi}) \end{aligned} \quad (\text{B.6})$$

where for sufficiently small ϵ_2 this equality holds for some $\underline{\phi}_d \in (\underline{\phi}^* - \epsilon_2, \underline{\phi}^*)$. Now since

$$\int_{\underline{\phi}^* - \epsilon_2}^{\underline{\phi}^* - \epsilon_1} \frac{1}{\delta} d\underline{\phi} = \ln \frac{\epsilon_2}{\epsilon_1} \quad (\text{B.7})$$

obviously the integral in (B.6) does not converge as $\epsilon_1 \rightarrow 0$. This implies, that if the factor $p(\underline{\theta}|\underline{\phi}) = |\underline{\Xi}'\underline{\Xi}|^{1/2}$ is left out of the prior, i.e. if prior independence between $\underline{\theta}$ and $\underline{\phi}$ is assumed, then the resulting posterior distribution $p_1(\underline{\phi}|\underline{y})$ concentrates all posterior probability in the neighbourhood(s) of the singularity point(s) if any, regardless of the data. In particular models containing a mean have such a singularity point namely $\underline{\phi}^* = 1$.

This clearly unacceptable consequence of adopting $p(\underline{\theta}|\underline{\phi}) = p(\underline{\theta})$ as a prior for $\underline{\theta}$ in the general model, of course also arises in the context of the stationary model.

Probing for Serial Correlation in Least Squares Regression.

Writing the regression model as

$$\begin{matrix} \tilde{y} &= & X & \theta & + & \tilde{e} & , \\ n \times 1 & & n \times p & p \times 1 & & n \times 1 \end{matrix} \quad (3.1)$$

where \tilde{y} is the vector of observations, X is the matrix of independent variables, and θ is the vector of the p linear parameters; it is usually assumed, that the random errors \tilde{e}_i , $i = 1, 2, \dots, n$, are independently distributed as $N(0, \sigma^2)$, and this assumption provides justification for the "standard" techniques of analysis.

For data collected in time (or space) sequence it may often be true, that the errors are in fact serially dependent, in particular this would seem likely, for many economic and business series.

Testing for serial correlation in situations of this kind, has been given much attention by a number of authors, most notably by Durbin and Watson [1950], [1951], [1971], whose well known test (DW-test for short) enjoys widespread use in applied work in many fields. In testing the null hypothesis of independence, the DW-statistic is constructed with a first order autoregressive (AR-1) noise structure (3.2) as the alternative hypothesis in mind.

$$\tilde{e}_i = \phi \tilde{e}_{i-1} + \epsilon_i \quad (3.2)$$

where the ϵ_i 's are i.i.d. $N(0, \sigma^2)$, i.e. they are a white noise series (also independent of \tilde{x}).

Modifications and new tests have been proposed from time to time, for example by Abrahamse and Louter [1971], Berenblut and Webb [1973], Durbin [1969], [1970], Grady [1970] and Schmidt [1972]. From the point of view of power, however, none of these alternatives have been

demonstrated to have any decisive advantage over the DW-test, not even when the true noise structure is different from AR-1 and the competing test was devised for such situations, Smith [1976] (also Abrahamse and Louter [1971], Berenblut and Webb [1973], Durbin and Watson [1971]).

The whole idea of testing for serial correlation in the circumstances in which it is usually done is somewhat suspect. For sequential data such as is often obtained in economics for example, one would usually expect serial correlation. There is therefore no particular reason for believing the null hypothesis, and it would often seem much more natural to assume, that the noise vector \tilde{e} in (3.1) is generated by an AR-1 scheme like (3.2), rather than by a white noise process.

It is the topic of this chapter to see how a broader appreciation of the inferential situation concerning an autoregressive parameter ϕ may be achieved, from a likelihood as well as a Bayesian point of view. In particular when ϕ behaves as a nuisance parameter in relation to θ in (3.1), to consider how inferences about θ may be made taking ϕ into account.

In cases where the null hypothesis has meaning, the likelihood and the Bayesian approaches suggest alternative testing procedures, which can be viewed as competitors to the DW-test. A simulation study is undertaken to determine their relative powers in conjunction with two selected model structures, as well as to shed light on other aspects of the alternatives.

Three examples are presented to illustrate consequences of employing different approaches.

3.1 The Durbin-Watson test.

Durbin and Watson's test statistic is defined as

$$d = \frac{\sum_{i=1}^{n-1} (r_i - r_{i+1})^2}{\sum_{i=1}^n r_i^2} \quad (3.3)$$

where the r_i 's are the residuals after an ordinary least squares fit, i.e.

$$r = (I - X(X'X)^{-1}X')y \quad (3.4)$$

and it is assumed without loss of generality that X has full column rank, p . The statistic d takes values between 0 and 4. Values close to 2 imply that the e_i 's in the linear model (3.1) are independent, while positive serial correlation will tend to produce smaller values and negative serial correlation larger ones. Unfortunately the sampling distribution of d under the null hypothesis of independence depends in general on X , so no single significance point $d_S(\alpha)$ depending only on n and p can be given. But bounds have been found $d_L(\alpha) \leq d_S(\alpha) \leq d_U(\alpha)$. The one-sided DW-test for serial correlation at the α level, thus consists of referring d to these tabulated bounds. If $d < d_L(\alpha)$ significance has been established, i.e. the hypothesis is to be rejected at the level α . If $d_L(\alpha) < d < d_U(\alpha)$ the test is said to be inconclusive; and employing popular terminology, the null hypothesis is to be accepted when $d_U(\alpha) < d$, although no more than lack of significance has been established at the given level.

Testing for negative serial correlation is done by referring the statistic (4-d) to the same bounds, and a two tail test at the 2α level results from doing both. The test would ordinarily be credited with the power associated with $d_L(\alpha)$ only; an exception is polynomial regression

where it has been shown that d_U gives a very good approximation to the true significance points of d , Hannan [1957].

$$\text{Rewriting (3.3) slightly as}$$

$$d = \frac{\sum_{i=2}^n x_i^2 + \sum_{i=1}^{n-1} x_i^2 - 2 \sum_{i=1}^{n-1} x_i x_{i+1}}{\sum_{i=1}^n x_i^2} = 2 - 2 \frac{\sum_{i=1}^{n-1} x_i x_{i+1}}{\sum_{i=1}^n x_i^2} \quad (3.5)$$

it is recognized that the fraction on the far right hand side is of the same form as (2.49) and (2.44), and hence provides an estimate of $\hat{\phi}$ which is essentially equivalent to performing only one iteration in the scheme of Figure 2.1 in Chapter 2. Denoting this estimate by $\hat{\phi}$, i.e.

$$\hat{\phi} = 1 - \frac{d}{2} \quad (3.6)$$

the DW-test might be formulated in terms of $\hat{\phi}$, in which case the test consists of determining whether $\hat{\phi}$ is significantly different from zero. This formulation is perhaps somewhat preferable, because the corresponding bounds help to give a direct indication of how pronounced the serial correlation must be in order to be detected. For example if $n = 15$ and $p = 2$, we find $\hat{\phi}_L(.05) = 1 - d_U(.05)/2 = .46$, implying that autocorrelation caused by a true value of say .3, might very well pass unnoticed.

The DW-test presupposes, that the linear regression model under consideration includes a mean. Durbin and Watson argue [1950], that this is only a slight limitation since a model lacking a mean can be augmented to include one for the purpose of conducting the test, and subsequently this extra independent variable may be eliminated again. Cochrane and Orcutt [1949] have shown, that the residuals \hat{r} from an ordinary least squares fit are biased towards randomness, and increasingly so as the number of independent variables is expanded. This

result would lead one to infer, contrary to Durbin and Watson, that the evidence of serial correlation carried by \tilde{x} , would tend to be obscured by forcing a mean into the linear model. Such an effect is indeed demonstrated in Sections 3.5 and 3.6.

Durbin and Watson do not address the question of what to do if the null-hypothesis is rejected. Cochran and Orcutt [1949] demonstrated that the ill effects of ignoring positive serial correlation of the errors in least squares regression might be remedied by applying an independence inducing transformation, and analyzing the transformed linear model by "standard" methods. If the error follows an AR-1 scheme, independence is induced by multiplying (3.1) by the matrix c :

$$c\tilde{y} = cX\tilde{\theta} + c\tilde{\epsilon} \quad (3.7)$$

$$\tilde{z} = \tilde{\epsilon} \tilde{\theta} + \tilde{\epsilon} \quad (3.8)$$

$$(n-1) \times 1 \quad (n-1) \times p \quad p \times 1 \quad (n-1) \times 1$$

where

$$c = \begin{bmatrix} -\phi & 1 & 0 & 0 \\ 0 & 1 & -\phi & 0 \\ & & \ddots & \ddots \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (3.9)$$

Theil and Nagar [1961] suggest to proceed as if c was known and was given by substituting $\hat{\phi}$ for ϕ in (3.9) (actually they employ a slight variant of $\hat{\phi}$). They caution, that not only does the approximation $\hat{\phi} = \phi$ neglect sampling variation in $\hat{\phi}$; but also that the same data are used to test the null hypothesis of independence, to estimate $\hat{\phi}$ if this hypothesis is rejected, and to reestimate $\hat{\theta}$ in that case. Hence they conclude, that the usual standard deviations for the least squares estimates of the θ_j 's in (3.8) may, on the average, underestimate the unreliability of these estimates.

3.2 $SS(\hat{\theta}, \hat{\phi})$ as a function of $\hat{\phi}$.

Before turning to the likelihood and Bayesian treatment, it is important to study $SS(\hat{\theta}, \hat{\phi})$ as a function of $\hat{\phi}$, where $SS(\hat{\theta}, \hat{\phi})$ is the residual sum of squares from fitting by least squares the transformed regression model (3.8), i.e.

$$SS(\hat{\theta}, \hat{\phi}) = \tilde{z}'(I - E(E'E)^{-1}E')\tilde{z} \quad (3.10)$$

It is assumed without loss of generality that X has full column rank, p ; however $E = cX$ need not have full column rank for all $\hat{\phi}$. It is therefore of interest to see how $SS(\hat{\theta}, \hat{\phi})$ behaves when $\hat{\phi}$ passes through singularity points (if any).

3.2.1 Singularity due to a mean.

If the original linear model (3.1) involves a mean, say θ_1 , so that

$$X = \begin{bmatrix} 1 & X_{(1)} \\ \vdots & \vdots \\ 1 & X_{(1)} \end{bmatrix} \quad (3.11)$$

where $\mathbf{1}$ is a vector of ones, then

$$E = cX = \begin{bmatrix} 1-\phi & & & \\ 1-\phi & cX_{(1)} & & \\ \vdots & \vdots & \ddots & \\ 1-\phi & & & \end{bmatrix} = \begin{bmatrix} \delta & E_{(1)} \\ \vdots & \vdots \end{bmatrix} \quad (3.12)$$

where $\delta = \delta \mathbf{1} = (1-\phi)\mathbf{1}$.

Clearly the transformed model (3.8) has less than full rank at the point $\phi = 1$ ($\delta = 0$), and the matrix $E'E$ becomes singular. The question is how $SS(\hat{\theta}, \hat{\phi})$ behaves as $\hat{\phi} \rightarrow 1$.

Restricting attention to a neighbourhood around $\hat{\phi} = 1$, where $E_{(1)}$ has full rank equal to $p-1$ (which is a proper restriction as proven below), we first write

and we may write

$$\tilde{X}_j = \begin{bmatrix} x_{2,j} - \phi x_{1,j} \\ \vdots \\ x_{n,j} - \phi x_{n-1,j} \end{bmatrix} = \begin{bmatrix} (\phi - \phi)x_{1,j} \\ \vdots \\ (\phi - \phi)x_{n-1,j} \end{bmatrix} = (\phi - \phi)X_j = \phi X_j + (3.22)$$

where $\delta = \phi - \phi$.

The problem is now of exactly the same form as above in subsection

3.2.1. Thus the conclusion may be drawn as before, that $SS(\tilde{\phi}, \phi)$ has a limit for $\phi \rightarrow \phi$, so that $SS(\tilde{\phi}, \phi)$ becomes continuous through $\phi = \phi$ by defining

$$SS(\tilde{\phi}, \phi) = \lim_{\phi \rightarrow \phi} SS(\tilde{\phi}, \phi). \quad (3.23)$$

Again, provided \tilde{X}_j is linearly independent of $\tilde{X}_1, \dots, \tilde{X}_{j-1}, \tilde{X}_{j+1}, \dots, \tilde{X}_p$ when $\phi = \phi$, the limiting value (3.23) may be interpreted as a residual sum of squares from fitting the linear model (3.8) with $\phi = \phi$, where a new parameter $\tilde{\phi}_j$ has been substituted for the one that became redundant for $\phi = \phi$.

It is noted, that at most one column may vanish for a particular value ϕ of ϕ . It is seen from (3.21) that ϕ determines \tilde{X}_j up to a multiplication constant, so that the occurrence of more than one column vanishing for $\phi = \phi$ is contrary to the assumption of \tilde{X} having full column rank.

3.2.3 General singularity.

In general, singularity may occur at $\phi = \phi$ when

$$k_1 \tilde{X}_1 + k_2 \tilde{X}_2 + \dots + k_p \tilde{X}_p = 0. \quad (3.24)$$

Assuming without loss of generality that $k_1 = -1$ we have equivalently

$$\tilde{X}_1 + \phi \tilde{X}_1 = k_2 (\tilde{X}_2 + \phi \tilde{X}_2) + \dots + k_p (\tilde{X}_p + \phi \tilde{X}_p). \quad (3.25)$$

This expression shows, that a situation of singularity cannot persist over a range of ϕ values. (3.25) is equating two first order polynomials in ϕ ($n-1$ times over), and can only be satisfied throughout an interval if the $(n-1)$ sets of coefficients are identical, i.e.

$$\begin{cases} \tilde{X}_1 = k_2 \tilde{X}_2 + \dots + k_p \tilde{X}_p \\ \tilde{X}_1 = k_2 \tilde{X}_2 + \dots + k_p \tilde{X}_p \end{cases} \quad (3.26)$$

$$\tilde{X}_1 = k_2 \tilde{X}_2 + \dots + k_p \tilde{X}_p \quad (3.27)$$

which contradicts the assumption of $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_p$ being linearly independent.

If we define

$$\begin{cases} \tilde{X}_1 = \tilde{X}_1 - k_2 \tilde{X}_2 - \dots - k_p \tilde{X}_p \\ \tilde{X}_1 = \tilde{X}_1 - k_2 \tilde{X}_2 - \dots - k_p \tilde{X}_p \end{cases} \quad (3.28)$$

then (3.25) may be expressed as

$$\tilde{X}_1 = \phi \tilde{X}_1 \quad (3.29)$$

which (much like (3.21)) determines \tilde{X}_1 up to a multiplicative constant, where

$$\tilde{X}_1 = \tilde{X}_1 - k_2 \tilde{X}_2 - \dots - k_p \tilde{X}_p. \quad (3.30)$$

To see that a second or higher order singularity (i.e. $\text{rank}(\tilde{X}) < p-1$ for $\phi = \phi$) can never occur, let us suppose that a second linear decomposition of \tilde{X}_1 existed:

$$\tilde{X}_1 = k'_2 \tilde{X}_2 + \dots + k'_p \tilde{X}_p \quad (3.31)$$

then

$$\tilde{X}_1 = \tilde{X}_1 - k'_2 \tilde{X}_2 - \dots - k'_p \tilde{X}_p \quad (3.32)$$

would of course also satisfy (3.29), which implies that \tilde{X}_1 and \tilde{X}_1 are identical up to a multiplicative constant. But since $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_p$ are linearly independent this can only happen if $k_2 = k'_2, \dots, k_p = k'_p$.

Finally defining

$$\begin{aligned} \tilde{X} &= [\tilde{X}_1 \tilde{X}_2 \dots \tilde{X}_p] \\ \text{with } \tilde{X}_1 &= \tilde{X}_1 + \phi \tilde{X}_1 \end{aligned} \quad (3.33)$$

we have in a neighbourhood around $\phi = \phi$, that

$$\begin{aligned}
 SS(\hat{\theta}, \phi) &= \tilde{z}'(I - \tilde{E}(\tilde{E}'\tilde{E})^{-1}\tilde{E}')\tilde{z} \\
 &= \tilde{z}'(I - \tilde{E}'\tilde{E}(\tilde{E}'\tilde{E})^{-1}\tilde{E}')\tilde{z}
 \end{aligned}
 \quad (3.34)$$

hence the general singularity of \tilde{E} at $\phi = \phi^*$ is now expressed in terms of a vanishing column \tilde{E}_1 in \tilde{E}^* , i.e. as a simple singularity considered above.

3.2.4 Singularity in polynomial regression.

A situation of particular interest arises, when a mean θ_1 creates a singularity when $\phi = 1$, but the presence of a deterministic linear trend in the model also makes $\tilde{E}^{-1} = [1 \ \tilde{E}(1)]$ singular when $\phi = 1$, so that the last step in (3.16) does not go through. Still in a neighbourhood around $\phi = 1$ we have

$$\begin{aligned}
 SS(\hat{\theta}, \phi) &= \tilde{z}'(I - \tilde{E}(\tilde{E}'\tilde{E})^{-1}\tilde{E}')\tilde{z} \\
 &= \tilde{z}'(I - \tilde{E}'\tilde{E}(\tilde{E}'\tilde{E})^{-1}\tilde{E}')\tilde{z}.
 \end{aligned}
 \quad (3.35)$$

If without loss of generality we say that $\tilde{x}_2 = \begin{bmatrix} 1 \\ 2 \\ \vdots \\ n \end{bmatrix}$, then

$$\tilde{E}_2 = \begin{bmatrix} 1 + (1-\phi) \\ 1 + 2(1-\phi) \\ \vdots \\ 1 + (n-1)(1-\phi) \end{bmatrix}
 \quad (3.36)$$

now defining

$$\tilde{E}_2^* = [\tilde{E}_2^* \ \tilde{E}_3^* \ \dots \ \tilde{E}_p^*]
 \quad (3.37)$$

where

$$\tilde{E}_2^* = \tilde{E}_2^{-1} = \begin{bmatrix} (1-\phi) \\ 2(1-\phi) \\ \vdots \\ (n-1)(1-\phi) \end{bmatrix}
 \quad (3.38)$$

it is recognized, that letting θ_2 retain its interpretation as a

deterministic trend at the point $\phi = 1$, i.e. manifesting itself as a mean for the differences, the "replacement parameter" θ_1^* for θ_1 has the interpretation of a deterministic linear trend for the differences.

Higher order polynomial trends may obviously be handled similarly.

3.2.5 Summary of findings.

Lemma 1

$$\text{Let } \tilde{E} = \begin{bmatrix} \tilde{E}_1 & \tilde{E}_2 & \dots & \tilde{E}_p \\ (n-1) \times p & (n-1) \times n & n \times p & n \times 1 \end{bmatrix}$$

and let \tilde{X} be any real matrix of full column rank, p . Then \tilde{E} has also rank p except possibly at distinct points $\phi = \phi^*$, where $\text{rank}(\tilde{E}) = p-1$.

Lemma 2

For $\phi \neq \phi^*$, denote by $SS(\hat{\theta}, \phi)$ the residual sum of squares from fitting by least squares the (transformed) linear model $\tilde{z} = \tilde{E}\hat{\theta} + \tilde{\epsilon}$, where $\tilde{z} = \begin{bmatrix} \tilde{z}_1 \\ \tilde{z}_2 \\ \tilde{z}_3 \\ \tilde{z}_4 \end{bmatrix}$ is any n -dimensional vector (of "original" observations) and $\tilde{\epsilon}$ is the vector of random disturbances. Then $SS(\hat{\theta}, \phi)$ has a limiting value as $\phi \rightarrow \phi^*$, and it becomes a continuous function of $\phi \in R$ if, by definition, we set $SS(\hat{\theta}, \phi^*) = \lim_{\phi \rightarrow \phi^*} SS(\hat{\theta}, \phi)$.

Corollary

If $\tilde{E}_j = 0$ for $\phi = \phi^*$ then $SS(\hat{\theta}, \phi^*)$ may be computed as the residual sum of squares from fitting by least squares the linear model

$$\tilde{z} = \tilde{E}^* \hat{\theta} + \tilde{\epsilon} \text{ provided } \text{rank}(\tilde{E}^*) = p, \text{ where}$$

$$\tilde{E}^* = [\tilde{E}_1 \dots \tilde{E}_{j-1} \tilde{E}_{j+1} \dots \tilde{E}_p]$$

and \tilde{X}_j is the j -th column of \tilde{X} having deleted the n -th row. The interpretations of the θ 's are unchanged when $\phi = \phi^*$ except for θ_j

which at that point assumes a different role.

In particular if $\tilde{x}_j = 1$, i.e. the linear model involves a mean, then $\phi = 1$ is a singularity point. If further the model does not involve a deterministic trend, then at this point θ_1 has the interpretation of a mean for the differenced observations, $y_2 - y_1, y_3 - y_2, \dots, y_n - y_{n-1}$. If the model includes polynomial trends up to and including the order q , then θ_j has the interpretation of a q -th order polynomial trend for the differences.

3.3 Inference about ϕ based on maximized likelihood.

The likelihood function corresponding to the model (3.1) and (3.2)

may be written as (Equation (2.29) in Chapter 2):

$$l(M, \theta, \phi, \sigma | \tilde{y}) = \sigma^{-n} \exp \left(-\frac{1}{2} \sigma^{-2} \sum_{j=1}^n (z - \tilde{z})^2 \right) + (e_1 - M)^2 \quad (3.39)$$

where M is a starting parameter for the noise:

$$e_1 = M + \epsilon_1 \quad (3.40)$$

If it is known a priori that the noise is stationary and reversible, then $e_1 = (1 - \phi)^{-1/2} \epsilon_1$ (see for example Anderson [1954]); but since such knowledge will seldom be available, M is included to hedge against the possibility that the first observation might have an untypical error, and more importantly to broaden the noise structure to cover explosive as well as stationary AR-1 processes (see Section 2.1 of Chapter 2).

Maximizing the likelihood (3.39) with respect to $(M, \theta, \phi, \sigma)$ it is found that

$$\max_{(M, \theta, \phi, \sigma)} l(M, \theta, \phi, \sigma | \tilde{y}) = \left(\frac{SS(\hat{\theta}, \hat{\phi})}{n} \right)^{-n/2} e^{-n/2} \quad (3.41)$$

where

$$SS(\hat{\theta}, \hat{\phi}) = \min_{\theta} SS(\hat{\theta}, \hat{\phi}) \quad (3.42)$$

It is noted, that the maximum likelihood estimate (MLE) of ϕ , is the value $\hat{\phi}$ which minimizes the residual sum of squares left by an ordinary least squares fit of the transformed linear model (3.8).

A test for a hypothesis of $\phi = \phi^*$ may be constructed from the likelihood ratio

$$\lambda_{\phi^*} = \frac{\max_{(M, \theta, \phi^*, \sigma)} l(M, \theta, \phi^*, \sigma | \tilde{y})}{\max_{(M, \theta, \phi, \sigma)} l(M, \theta, \phi, \sigma | \tilde{y})} = \frac{l(\hat{M}_0, \hat{\theta}_0, \phi^*, \hat{\sigma}_0 | \tilde{y})}{l(\hat{M}, \hat{\theta}, \hat{\phi}, \hat{\sigma} | \tilde{y})} = \left(\frac{SS(\hat{\theta}_0, \phi^*)}{SS(\hat{\theta}, \hat{\phi})} \right)^{-n/2} \quad (3.43)$$

by considering the well known large sample χ^2 distribution approximation of the likelihood ratio statistic

$$\Lambda_{\phi^*} = -2 \ln \lambda_{\phi^*} \xrightarrow{d} \chi^2_1 \quad (3.44)$$

This approximation is justified as the asymptotic distribution of Λ_{ϕ^*} (under $\phi = \phi^*$), when the maximum likelihood estimators are asymptotically normally distributed (see e.g. Rao [1965] or Wilks [1962]). Usually the asymptotic normality of MLE is established for independent observations, but Whittle [1953] has extended the limiting properties of MLE to cover stationary time series. More specifically Hildreth [1969] showed, that the MLE of the parameters of a linear model with stationary first order autoregressive errors have an asymptotic multivariate normal distribution with mean vector equal to the true parameter values, and the estimators are asymptotically efficient. Hence the χ^2 approximation for the distribution of Λ_{ϕ^*} holds for $-1 < \phi^* < 1$.

For $|\phi^*| > 1$, i.e. for explosive series, it does not make sense to talk about asymptotic results, and in fact it is demonstrated in Section 3.6, that (3.44) does not hold in this case. Exceptions must also be made for singularity points $\phi^* = \phi$. For although $SS(\hat{\theta}, \hat{\phi})$ can be made continuous at these points, it was seen in Section 3.2 that

the continuity required an artificial inclusion of a substitution parameter, which is not part of the original model, hence a test for $\phi^* = \hat{\phi}$ using (3.44) would be conditional on an implicit misspecification of the model at this very point.

In particular (3.44) suggests an approximate likelihood ratio test for $\phi = 0$ based on

$$\Lambda_0 = -2 \ln \lambda_0 = n \ln SS(\hat{\theta}_0, \phi = 0) - n \ln SS(\hat{\theta}, \hat{\phi}) \sim \chi^2_1 \quad (3.45)$$

(How well this approximation works in moderate sample sizes is studied by simulation in Section 3.6).

In practice rather than merely testing whether Λ_0 takes a significantly large value it is usually most informative to plot

$$LMAX(\hat{\phi}) = -n \ln SS(\hat{\theta}, \hat{\phi}) \quad (3.46)$$

over a suitable range of $\hat{\phi}$. The mode of the $LMAX(\hat{\phi})$ -curve marks $\hat{\phi}$ and an approximate $100(1-\alpha)\%$ confidence interval for $\hat{\phi}$ is obtained from

$$LMAX(\hat{\phi}) - LMAX(\hat{\phi}) < \chi^2_1(\alpha). \quad (3.47)$$

In other words the plausibility of any value of $\hat{\phi}$ may be approximately assessed from such a plot, and it becomes visually clear whether the data set at hand contains much information about $\hat{\phi}$ or only relatively little.

Ordinarily $\hat{\phi}$ is a nuisance parameter in relation to the linear parameters $\hat{\theta}$, and inferences about $\hat{\phi}$ may only be of interest as a stepping stone to the further analysis concerning $\hat{\theta}$. Unless strong reasons exist for believing that $\hat{\phi} = \phi^*(=0)$, the occurrence of $\hat{\phi}^*$ inside a certain confidence interval supplies little justification for outright assuming that $\hat{\phi} = \phi^*$. From a likelihood point of view it would seem more natural to base inferences about $\hat{\theta}$ on $\hat{\phi} = \hat{\phi}$, and

hence to estimate the independence inducing transformation c simultaneously with $\hat{\theta}$.

3.4 Bayesian inference about $\hat{\phi}$.

From the Bayesian point of view all that can be said about $\hat{\phi}$ in light of the data, is embodied in the marginal posterior distribution of $\hat{\phi}$, $p(\hat{\phi}|\underline{y})$. As was shown in Chapter 2, the posterior density function of $\hat{\phi}$ resulting from multiplying the likelihood function l in (3.39) by a prior $p(M, \hat{\theta}, \hat{\phi}, \sigma)$ which is approximately noninformative for $(M, \hat{\theta}, \sigma)$, and integrating out $(M, \hat{\theta}, \sigma)$ is

$$p(\hat{\phi}|\underline{y}) = p(\hat{\phi})(SS(\hat{\theta}, \hat{\phi}))^{-\frac{n-p-1}{2}} \quad (3.48)$$

where $p(\hat{\phi})$ is the marginal prior for $\hat{\phi}$.

If $p(\hat{\phi})$ is chosen uniform, then

$$p_u^e(\hat{\phi}|\underline{y}) = SS(\hat{\theta}, \hat{\phi})^{-\frac{n-p-1}{2}} \quad (3.49)$$

where "u" stands for "uniform" and "e" stands for "exact". Obviously $p_u^e(\hat{\phi}|\underline{y})$ has its mode at the least squares value $\hat{\phi} = \hat{\phi}$; and this distribution may be approximated closely by a t-distribution with $v = n-p-2$ degrees of freedom provided $SS(\hat{\theta}, \hat{\phi})$ is approximated near $\hat{\phi}$ by

$$Q_t(\hat{\phi}) = b_1 + b_2 \hat{\phi} + b_3 \hat{\phi}^2 \quad (3.50)$$

then relating to the t-form

$$p_u^e(\hat{\phi}|\underline{y}) \approx \left(1 + \frac{(\hat{\phi} - \hat{\phi})^2}{v s_t^2}\right)^{-\frac{v+1}{2}} = p_u(\hat{\phi}|\underline{y}) \quad (3.51)$$

we have that

$$p_u(\hat{\phi}|\underline{y}) \sim t(\hat{\phi}_t, s_t, v) \\ \text{with } \hat{\phi} = -\frac{b_2}{2b_3} \\ \theta_t = \frac{b_1 b_3 - b_2^2}{4b_3^2 v} \quad (3.51)$$

and $v = n-p-2$.

Alternatively in the vicinity of $\hat{\phi}$, an approximately noninformative prior for ϕ is (see Section 2.5 of Chapter 2)

$$P_n(\phi) \propto (n-1) + (n-2)\phi^2 + \dots + \phi^{2(n-2)} \quad (3.52)$$

where " ∞ " stands for "noninformative" (but also serves to remind, that this prior is a function of n). However a direct multiplication by

$P_n(\phi)$ makes the right hand side of (3.48) diverge as ϕ becomes large.

As discussed in Section 2.5 of Chapter 2, the argument leading to this approximately noninformative prior is a local one, applicable in the vicinity of $\hat{\phi}$ only. Here the effect on the posterior of the noninformative prior as compared to a uniform prior, is to shift its mode moderately

away from the origin, along with an increase of the spread. These modifications may be accomplished within the t-form by fitting the polynomial

$$Q_t = \beta_1 + \beta_2\phi + \beta_3\phi^2 \quad (3.53)$$

through the three points $Q_t(\phi' - \Delta\phi)$, $Q_t(\hat{\phi})$ and $Q_t(\phi' + \Delta\phi)$, calculated from

$$Q_t(\phi' \pm \Delta\phi) = q(\phi' \pm \Delta\phi) \cdot Q_t(\hat{\phi} \pm \Delta\phi) \quad (3.54)$$

where ϕ' is reasonably close to $\hat{\phi}$ and $\Delta\phi$ is about or less than s_t , and where

$$q(\phi) \propto P_n(\phi) \cdot \frac{2}{n-p-1} \quad (3.55)$$

The modified quadratic curve (3.53) determines

$$P_n(\phi|y) \sim t(\hat{\phi}, s_t, v) \quad (3.56)$$

where $\hat{\phi}$, s_t and v are as in (3.51) except the b 's are replaced by b 's.

The Bayesian parallel to a significance test at the α level of a hypothesis $\phi = \phi^*$, consists of determining whether ϕ^* is or is not

included in the $100(1 - \alpha)\%$ highest posterior density (HPD) interval for ϕ (see for example Box and Tiao [1972], pp. 122-127). This may now conveniently be checked by referring the quantities

$$t_{\phi^*} = \frac{\phi^* - \hat{\phi}}{s_t} \quad (3.57)$$

or

$$t_{\phi^*} = \frac{\phi^* - \hat{\phi}}{s_t} \quad (3.58)$$

to a $t(0, 1, n-p-2)$ distribution depending on which prior is adopted.

In this treatment the plausibility of any hypothesized value $\phi^* < \phi^0 < \infty$ may be "tested", except singularity points $\phi^0 = \hat{\phi}$. In particular this excludes (as before) $\phi^0 = 1$ for models involving a mean. This is unfortunate since $\phi^0 = 1$ implying that the noise is a random walk may be of special interest (this question shall be returned to in Chapter 4).

From a Bayesian viewpoint in making inferences about the linear parameters θ , the whole posterior distribution $p(\theta|y)$ should be used; and we have seen that utilizing the familiar t form, $p(\phi|y)$ may be approximated with a modest computational effort.

Specifically

$$p(\theta|y) = \int_R p(\theta|\phi, y) p(\phi|y) d\phi \quad (3.59)$$

where

$$p(\theta|\phi, y) \sim t(\hat{\theta}, s^2(\hat{\theta}|\phi)^{-1}, v)$$

with

$$\hat{\theta} = (E'E)^{-1}E'y$$

$$v s^2 = SS(\hat{\theta}, \phi)$$

and

$$v = n-p-1.$$

When the posterior (3.59) may be approximated by

$$p(\theta|y) \approx p(\hat{\theta}, y) \quad (3.61)$$

the Bayesian inference about θ exactly parallels the likelihood solution suggested in the preceding Section 3.3.

3.5 Examples.

In this section we shall reexamine three data sets used as examples in previously published papers.

The first set is one generated by Zellner and Tiao [1964], and the "true" process does not include a mean. The second set is the textile data examined in the paper by Theil and Nagar [1961]. The third set is the "spirits series" considered by Durbin and Watson [1951], and re-examined by Theil and Nagar [1961].

3.5.1 The data generated by Zellner and Tiao.

Zellner and Tiao [1964] generated and analyzed the data in Table 3.1, using the model:

$$Y_i = \theta x_i + e_i \quad \text{with} \quad \theta = 3. \quad (3.62)$$

$$\begin{cases} e_1 = M + \epsilon_1 & \text{with} \quad M = .5 \\ e_i = \phi e_{i-1} + \epsilon_i & \text{with} \quad \phi = .5 \end{cases} \quad (3.63)$$

where the ϵ_i 's are i.i.d. $N(0,1)$.

Probing for serial correlation along the lines of Section 3.3, $LMAX(\phi)$ (3.46) is plotted in Figure 3.1a. $LMAX(\phi)$ has its maximum at $\phi = .36$, and the upper horizontal bars in the plot marks the level of $LMAX(\phi)$. The lower set of bars are drawn in a distance of 3.84 = $\chi^2_{.05}(1)$ from $LMAX(\phi)$, hence an approximate 95% confidence interval for ϕ is read off as $-.14 < \phi < .88$. The evidence of positive serial correlation is not very strong, in particular Λ_0 is not significant at the 5% level. (The points $LMAX(\phi=0)$ and $LMAX(\phi=1)$ are indicated by vertical bars on the $LMAX(\phi)$ graph).

TABLE 3.1, Zellner and Tiao's data.

i	y_i	x_i
1	12.649	3.9
2	18.794	6.0
3	12.198	4.2
4	14.372	5.2
5	13.909	4.7
6	14.556	5.1
7	14.700	4.5
8	18.281	6.0
9	13.890	3.9
10	10.318	4.1
11	5.473	2.2
12	4.044	1.7
13	6.361	2.7
14	7.036	3.3
15	13.368	4.8

Figure 3.1b gives the Bayesian picture of what the data has to tell about ϕ . The solid curve is $p_n(\phi|y)$ with mode at $\hat{\phi} = .39$, and its 95% HPD interval for ϕ is $-.25 < \phi < 1.01$. The two broken curves in Figure 3.1b, $p^e_u(\phi|y)$ and $p^e_l(\phi|y)$, agree almost perfectly; also it makes little difference if inference about ϕ were based on those distributions rather than on $p_n(\phi|y)$.

Turning to the DW-test it is perhaps somewhat surprising to find $d = 2.002$, i.e. $\hat{\phi} = -.001$. Of course in order to compute d , it is necessary to artificially include a mean in the linear model, i.e. relate the data to the model

$$y_i = \theta_1 + \theta_2 x_i + e_i \quad (3.64)$$

where e_i is as in (3.63).

It is interesting to redraw $LMAX(\phi)$ and $p(\phi|\tilde{y})$ for this modified model (3.64); this is done in Figure 3.2. It discloses, that incorporating a mean into the model, destroys all traces of positive serial correlation which are imprinted in the data. Further insight into the problems created by model misspecification may be gained by considering what inferences can be drawn about the linear parameters θ_1 and θ_2 of (3.64). Figure 3.3a shows posterior distributions of the mean θ_1 . The graph $p(\theta_1|\phi=0, \tilde{y})$ is also the confidence distribution of θ_1 from a sampling theory point of view, since it is identical to the distribution of the statistic $\hat{\theta}_1$ conditional on $\phi = 0$. Since in this framework according to the procedure recommended by Durbin and Watson the value $\phi = 0$ would be accepted, then θ_1 would subsequently be found significantly different from zero at the $\alpha = .005$ level (two sided), and θ_1 could be mistakenly retained in the model. The distribution $p(\theta_1|\phi = \hat{\phi}, \tilde{y})$ can also be interpreted as a confidence distribution for θ_1 given that $\phi = \hat{\phi} = .36$. In this case the significance of θ_1 has dropped to the $\alpha = .03$ level (two sided). In contrast a Bayesian analysis does not discount the possibility of θ_1 being zero, as this point lies well within the 90% HPD region of $p_n(\theta_1|\tilde{y})$ (and of $p_0(\theta_1|\tilde{y})$).

It is of further interest to see what impact the inclusion of the mean, θ_1 , has on making inferences about the regression coefficient which appears as θ in the model (3.62) and as θ_2 in the model (3.64), and which is most likely to be the parameter of primary interest. Figures 3.3b and 3.3c show posterior distributions for θ_2 and θ respectively. It is seen, that not only has the inclusion of θ_1 the effect of diluting the information in the data about θ (θ_2), but $p(\theta_2|\phi = 0, \tilde{y})$ in Figure 3.3b actually throws doubt on the plausibility

MAXIMIZED LOG LIKELIHOOD

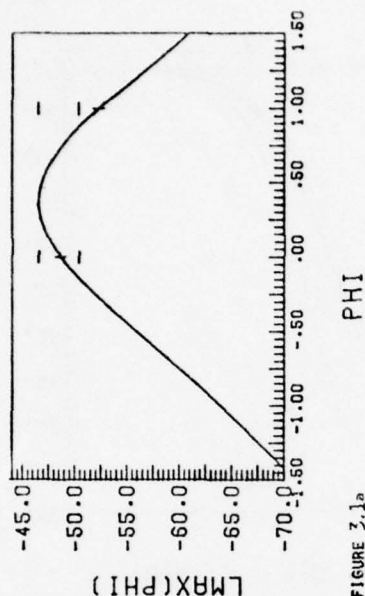


FIGURE 3.1a

MARG. POST. DIST. OF PHI

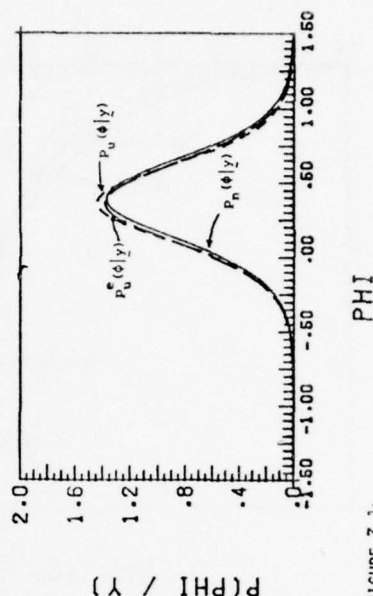


FIGURE 3.1b

MAXIMIZED LOG LIKELIHOOD

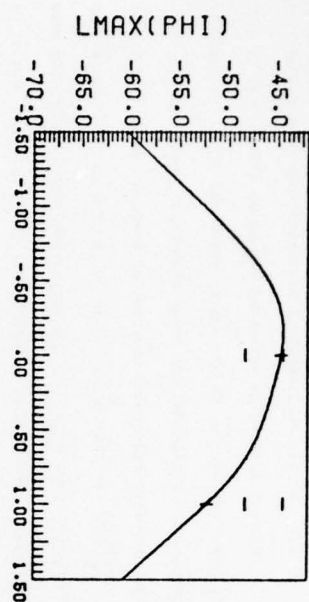


FIGURE 3.2a

PHI

MARG. POST. DIST. OF PHI

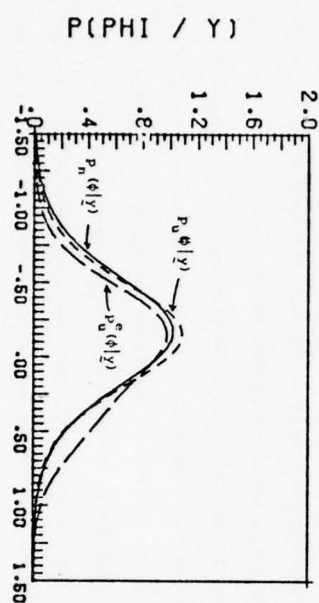


FIGURE 3.2b

PHI

MARG. POST. DIST. OF THETA

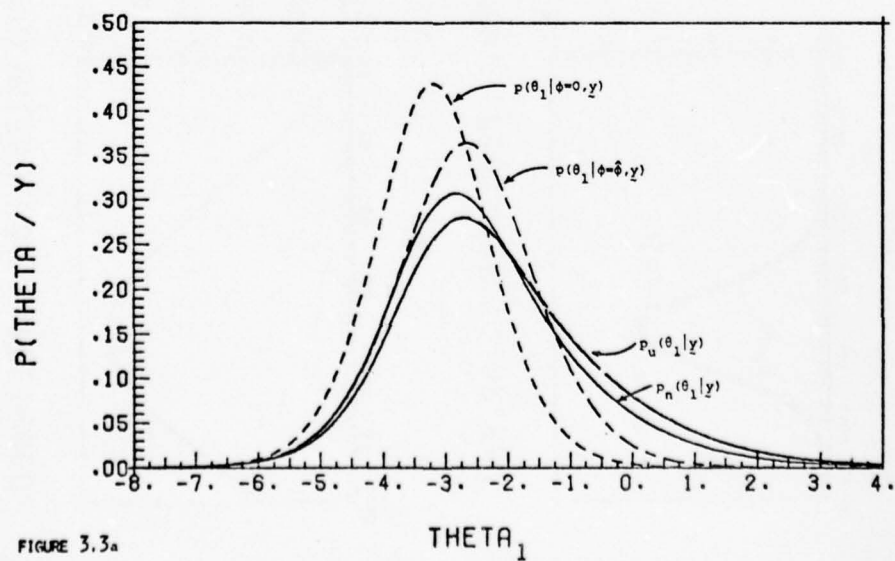


FIGURE 3.3a

THETA₁

MARG. POST. DIST. OF THETA

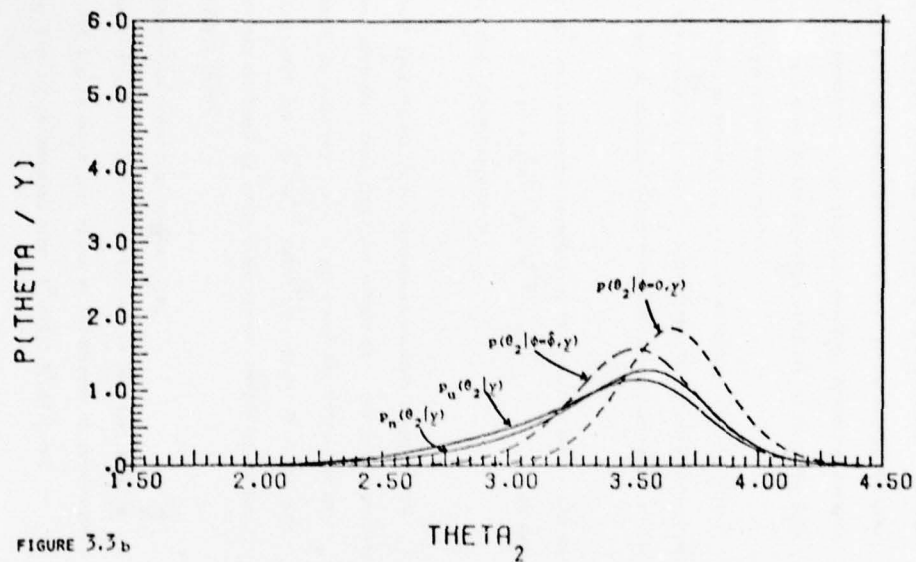


FIGURE 3.3b

97

MARG. POST. DIST. OF THETA

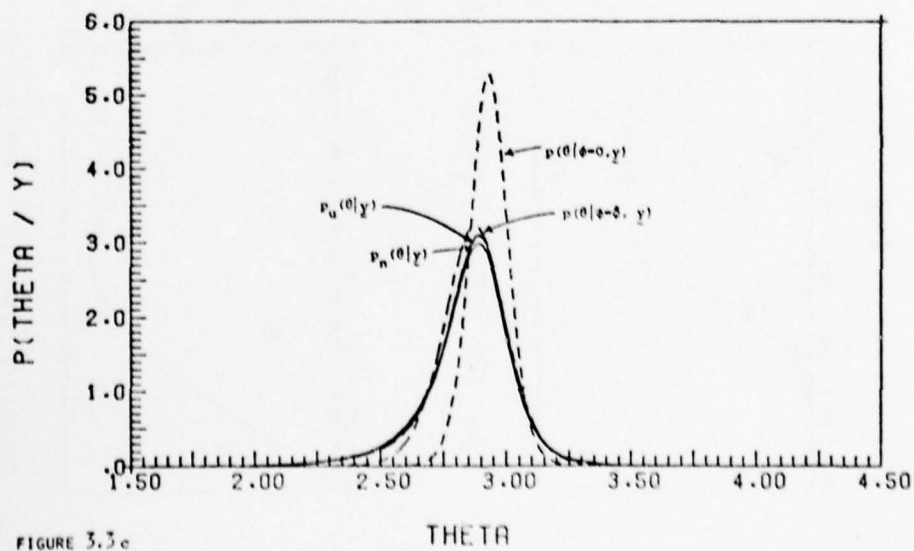


FIGURE 3.3c

98

of the true value $\theta_2 = 3$. On the other hand $p(\theta_2|y)$ covers $\theta_2 = 3$ nicely, while this value is just included in the 95% HPD region of $p(\theta_2|\phi = \hat{\phi}, y)$.

In Figure 3.3c it is observed that $p_n(\theta|y)$, $p_u(\theta|y)$ and

$p(\theta|\phi = \hat{\phi}, y)$ are almost coincidental, while an assumption of independence expressed through $p(\theta|\phi = 0, y)$ would lead one to believe, that θ is known with greater precision than is justified.

3.5.2 The textile data.

This data set consists of three simultaneous series of length 17 (yearly observations), viz. y_i , $x_{i,1}$ and $x_{i,2}$, where $y_i = \log$ (per capita consumption of textile), $x_{i,1} = \log$ (real per capita income), and $x_{i,2} = \log$ (deflated price index for clothing). The data shall not be relisted here, they appear along with their sources in Theil and Nagar [1961].

The model under consideration is

$$y_i = \theta_1 + \theta_2 x_{i,1} + \theta_3 x_{i,2} + e_i \quad (3.65)$$

where θ_2 and θ_3 are sometimes referred to as elasticities in economics.

Fitting (3.65) by ordinary least squares, the DW-statistic is calculated to be $d = 1.926$ ($\hat{\phi} = .04$), so that data set has been used as an example where the null hypothesis $\hat{\phi} = 0$ is accepted, and inference about the θ 's may be drawn accordingly.

Studying Figure 3.4 as in the previous example, it is found both from maximized likelihood plot and from the marginal posterior density of ϕ , that while a prior hypothesis of $\phi = 0$ is certainly not dismissed, a very wide range of other hypotheses $\phi = \hat{\phi}^*$ cannot be dismissed either. Specifically Figure 3.4a yields the approximate 95%

MAXIMIZED LOG LIKELIHOOD

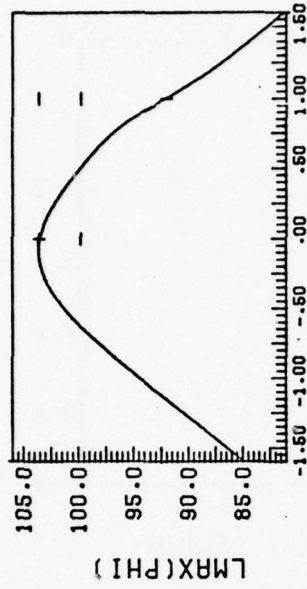


FIGURE 3.4a

MARG. POST. DIST. OF PHI

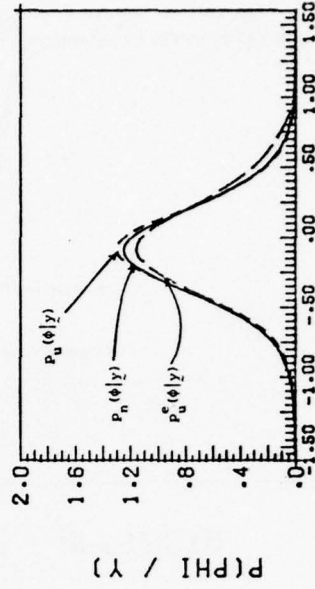


FIGURE 3.4b

POST. MARG. DIST. OF THETA

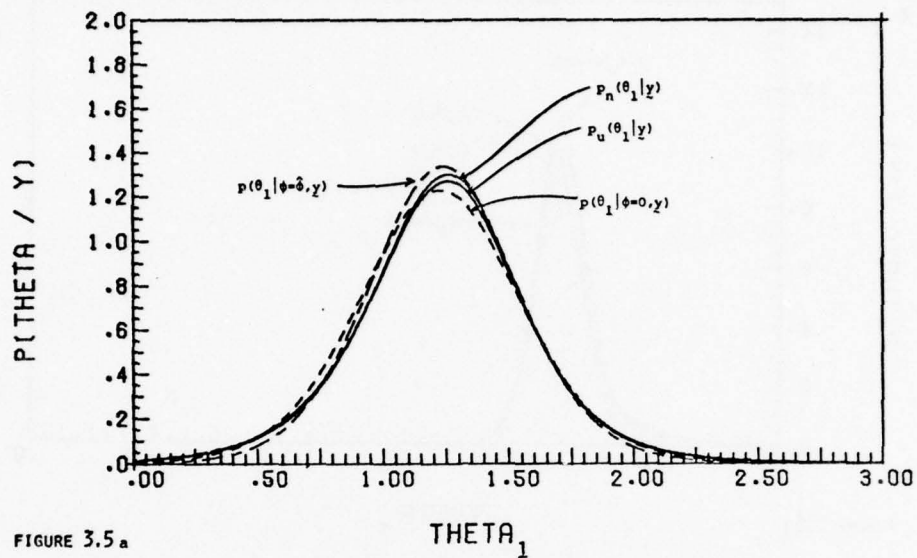


FIGURE 3.5a

THETA₁

POST. MARG. DIST. OF THETA

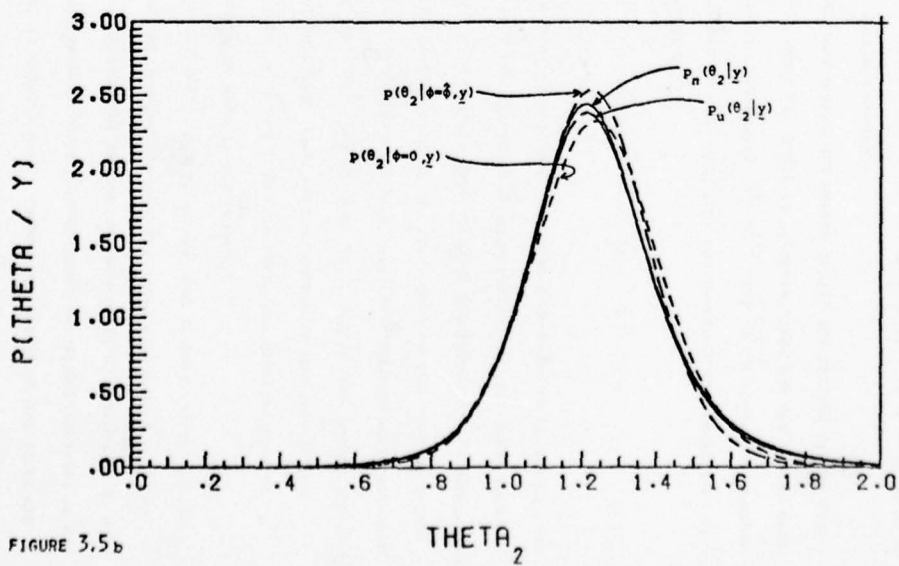


FIGURE 3.5b

THETA₂

confidence interval $-.63 < \hat{\phi} < .51$, with $\hat{\phi} = -.09$; and Figure 3.4b show a 95% HPD region of $-.78 < \hat{\phi} < .59$, with $\hat{\phi} = -.10$.

The question may be raised, whether there is any prior reason to suspect that ϕ is exactly zero. Causes can easily be put forth why observations of this kind would be positively serially correlated (e.g. a gradual change in purchasing habits) or negatively correlated (e.g. a cold winter may spur purchasing, with the after effect of a sluggish market the following year); but it is not easy to rationalize why the observations should be exactly independent.

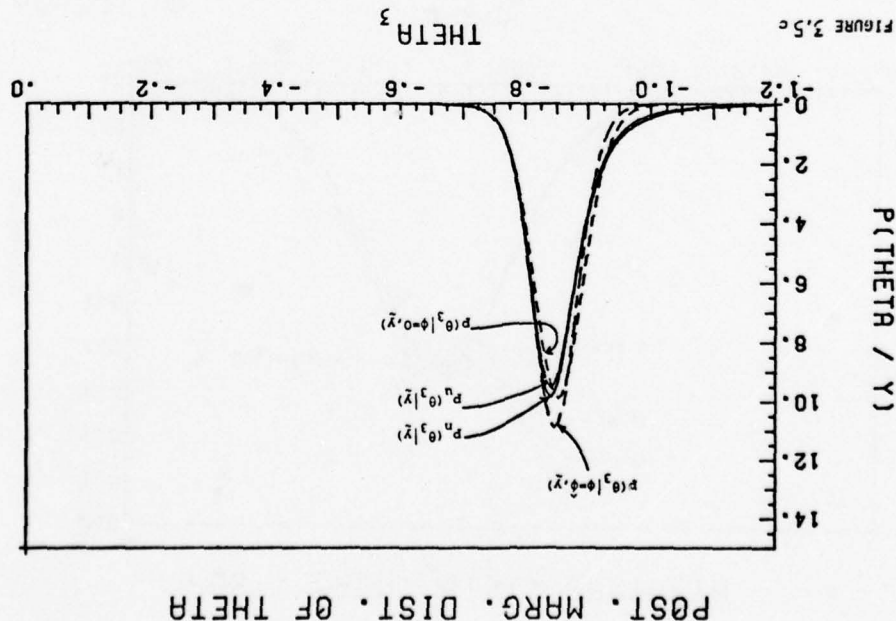
Figures 3.5a, 3.5b and 3.5c show posterior distributions of θ_1 , θ_2 and θ_3 respectively. For each θ_j parameter the four curves $p_n(\theta_j | y)$, $p_n(\theta_j | \tilde{y})$, $p(\theta_j | \hat{\phi} = \hat{\phi}, y)$ and $p(\theta_j | \hat{\phi} = 0, y)$ are almost identical. It is therefore a fact, that for this particular example where the mean of the posterior distribution of ϕ lies close to zero correct inferences about the θ_j 's result from assuming independence. It would seem however, that this has little to do with testing and accepting an unlikely hypothesis of $\phi = 0$, but is rather a consequence of the circumstance that for this data

$$p(\theta_j | \hat{\phi} = 0, y) \approx p(\theta_j | \hat{\phi} = \hat{\phi}, y) \approx p(\theta_j | y) \quad j = 1, 2, 3 \quad (3.66)$$

3.5.3 The spirits data.

This data set is very much like the textile data, consisting of three simultaneous time series y_1 , $x_{1,1}$ and $x_{1,2}$ defined as before except "spirits" takes the place of "textile", and now the series have 69 observations (see Durbin and Watson [1951] for listing of the data and reference to their source).

Fitting (3.65) to this data, the DW-statistic is calculated to be $d = .249$ ($\hat{\phi} = .875$), which is significant at the 2% level (two sided).



MAXIMIZED LOG LIKELIHOOD

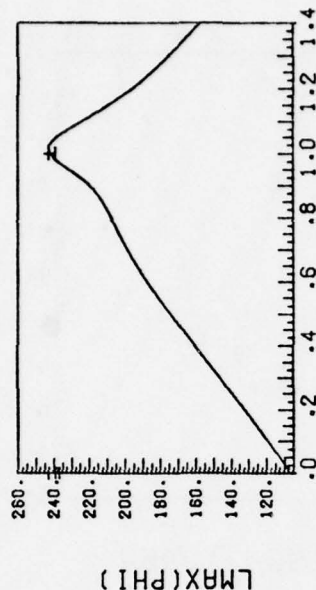


FIGURE 3.6a

MARG. POST. DIST. OF PHI

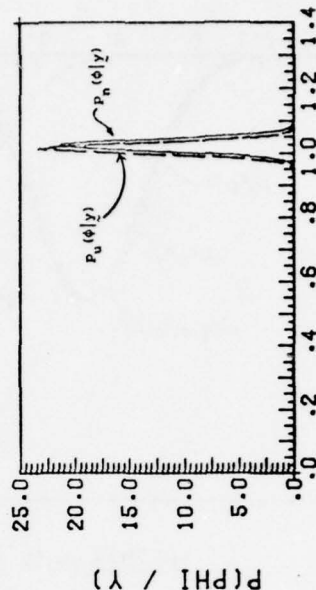


FIGURE 3.6b

105

The maximized likelihood curve in Figure 3.6a shows that $\hat{\phi} = 1.015$. An approximate 95% confidence interval constructed as described in Section 3.5.1 is $.98 < \phi < 1.05$; however this interval cannot be relied on, since the χ^2 approximation (3.44) only finds theoretical justification for $-1 < \phi < 1$. The Bayesian probe for serial correlation, Figure 3.6b, determines from $p_n(\phi|y)$ a 95% HPD interval of $.99 < \phi < 1.06$, with $\hat{\phi} = 1.025$.

Uniform agreement prevails that the observations are strongly serially correlated. Theil and Nagar [1961] carry out the further analysis concerning $\hat{\theta}$ as if $\hat{\phi} = \hat{\phi}$. But clearly $\hat{\phi}$ is a very poor estimate of ϕ . Approaching the problem from a sampling theory point of view it would seem most reasonable to use $\hat{\phi}$ rather than $\hat{\phi}$ in the independence inducing transformation (3.9) and then analyse (3.8) accordingly. Because while $\hat{\phi}$ is the maximum likelihood estimate of ϕ , $\hat{\phi}$ is a maximum likelihood estimate conditional on $\hat{\theta} = \hat{\theta} = (X'X)^{-1}X'y$ which may be a quite misleading estimate if ϕ is not close to zero. In the Bayesian framework when $p(\phi|y)$ is very narrow as here, then the conditional ($\phi = \hat{\phi}$, say) distribution of $\hat{\theta}$ will generally be quite similar to the marginal distribution.

The results of alternative approaches are shown in Figures 3.7a, 3.7b and 3.7c for θ_1 , θ_2 and θ_3 respectively. For θ_2 and θ_3 the three curves $p_n(\theta_j|y)$, $p_u(\theta_j|y)$ and $p(\theta_j|\hat{\phi}, y)$ are very much alike, where the latter one also may be interpreted as a confidence distribution for θ_j . An assumption of independence renders a grossly erroneous impression of the linear parameters $\hat{\theta}$; but also the substitution of $\hat{\phi} = \hat{\phi}$, while not quite as bad as $\hat{\phi} = 0$, yields misleading results. And it is not just a matter of overestimating the precision

MARG. POST. DIST. OF THETA

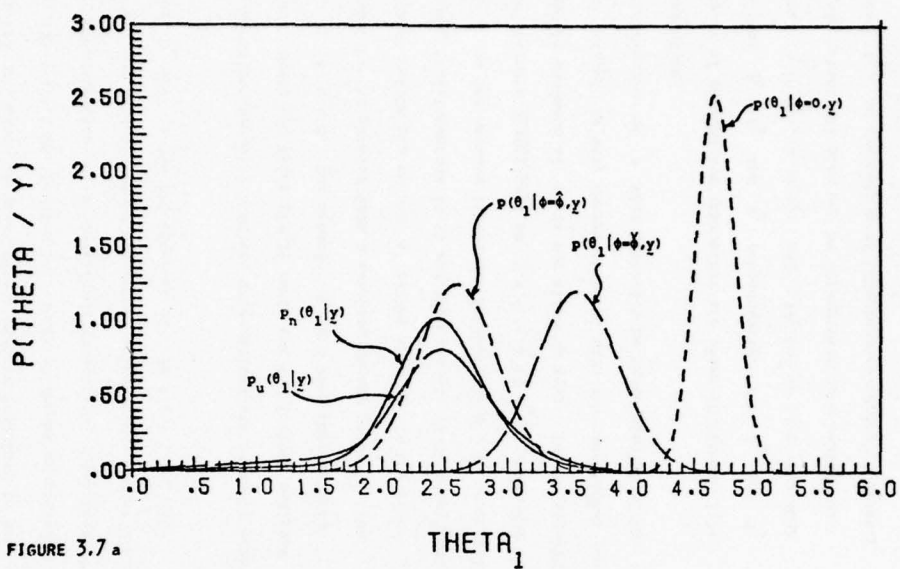


FIGURE 3.7 a

107

MARG. POST. DIST. OF THETA

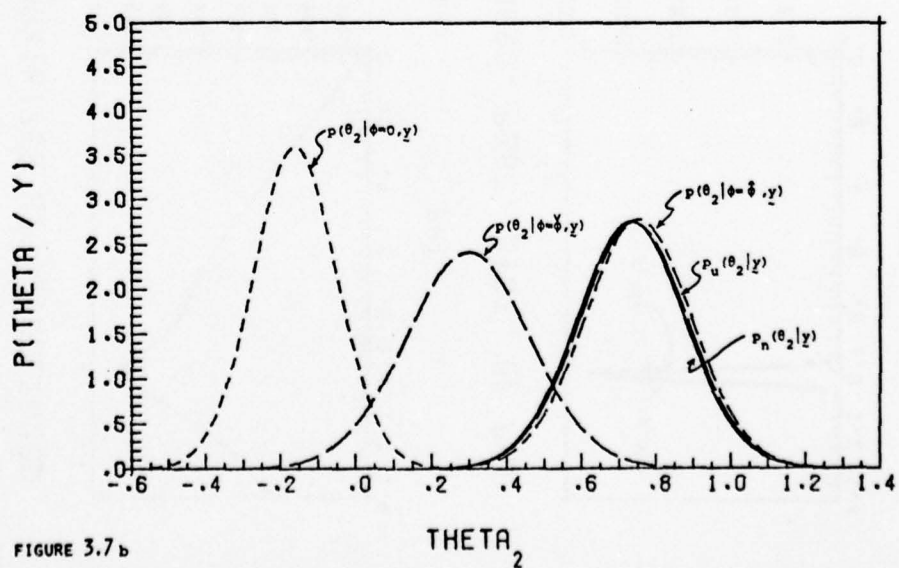


FIGURE 3.7 b

108

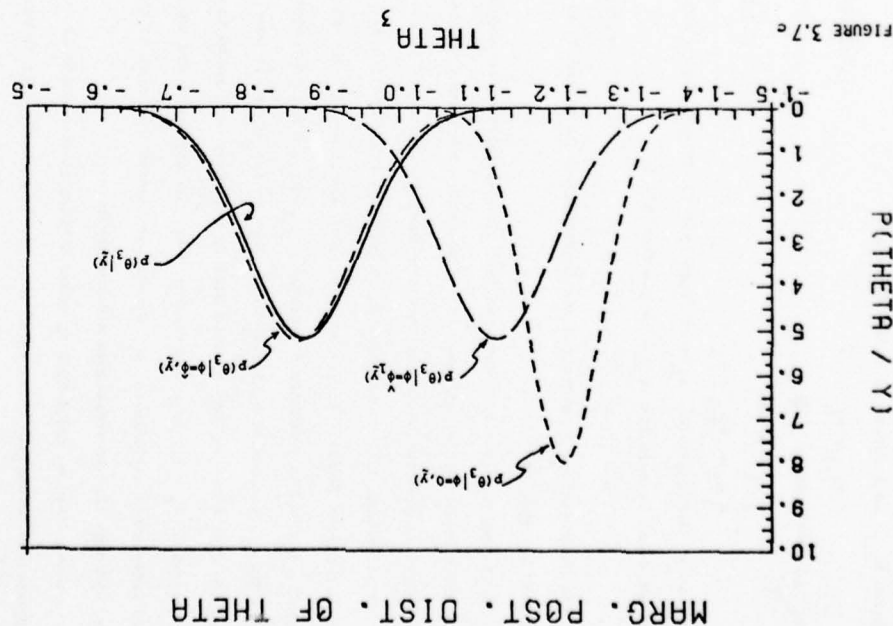


FIGURE 3.7

with which the θ_j 's are known, but one of location. Specifically the income elasticity is actually much greater than was suggested by Theil and Nagar's analysis, while the price elasticity has a somewhat smaller absolute magnitude.

Since the noise apparently is slightly explosive, the mean θ_1 is not a stable level (having allowed for the independent variables), but may more descriptively be termed an unstable equilibrium. The direction of the explosion is downwards, which may be interpreted as a development in time of consumer habits away from spirits (since logged data are analyzed the explosion is headed towards zero consumption).

Whether it would have been appropriate to difference the series before subjecting them to analysis is an unanswered question at this point. It shall be returned to in Chapter 4, and specifically we shall argue in the Appendix of that chapter, that differencing of this data is not desirable.

3.6 Monte Carlo study.

In this section the relative merits of the testing procedures associated with d , Λ_{ϕ^0} , t_{ϕ^0} and τ_{ϕ^0} are assessed by means of a simulation study. Although t_{ϕ^0} and τ_{ϕ^0} are derived at on Bayesian arguments, and therefore do not require any justification by sampling properties, we shall here study how they behave in repeated sampling.

The data generating models, A and B, employed in the present study are

$$A: Y_i = \theta_1 + \theta_2 x_{i,2} + e_i; \quad i = 1, 2, \dots, n \quad (3.67)$$

$$\text{with } \theta_1 = 10, \quad \theta_2 = 2.$$

and where $x_{i,2}$ is a random walk

$$x_{i,2} = \sum_{j=1}^i a_j \quad (3.68)$$

$$B: \quad Y_i = \theta x_{i,1} + e_i; \quad i = 1, 2, \dots, n \quad (3.69)$$

with $\theta = 2$.

and where $x_{i,1}$ is a series of random deviates

$$x_{i,1} = a_i \quad (3.70)$$

and e_i is AR-1 noise:

$$\begin{cases} e_1 = M + \epsilon_1 \\ e_i = e_{i-1} + \epsilon_i \end{cases} \quad (3.71)$$

The ϵ_i 's of (3.71) and the a_i 's of (3.68) and (3.70) are i.i.d. $N(0,1)$.

M was set equal to zero for $\phi \geq 1$, while for $\phi < 1$,

$$M = ((1-\phi)^{-2} - 1/2) \epsilon_1, \quad (3.72)$$

i.e. the noise follows a reversible stationary scheme ((2.5) and (2.6) in Chapter 2) when $\phi < 1$.

It is easy to visualize many practical circumstances where models like A and B would emerge, at least as a first step toward more comprehensive process descriptions. For instance Model A might describe economic relationships like those examined in Sections 3.5.2 and 3.5.3.

As Model A stands (3.67), θ_2 is an elasticity equal to 2, and the "shocks", a_i , in the independent variable $x_{i,2}$ have a variance equal to that of ϵ_i associated with the noise. But as far as ϕ is concerned the data generated by A could equally well have originated from a process with $\theta_2 = 1$, say, and $\text{Var}(a_i) = 4$. In any event the impact of the independent variable on y_i is more pronounced than that of the noise, although the two components are comparable.

A model like B might turn up under similar circumstances if it is appropriate to analyze differenced data, as is commonly done in econometrics.

From each model 400 independent samples were generated for sample sizes of $n = 15, 25, 60$ and for each $\phi = -.5, -.4, \dots, 1.1, 1.2$, except for

the combination $n = 60$ and $\phi = 1.2$, i.e. a total of $400 \times (2 \times 18 + 1 \times 17)$ independent samples for each model. 400 additional samples were generated by each model with $\phi = 0$ and for $n = 15, 25, 60$.

The models A and B generated data in parallel, i.e. they employed the same a_i 's and ϵ_i 's for corresponding samples (same ϕ , same n). Essentially all calculations were done in double precision.

3.6.1 Power curves.

Figures 3.8 (for $n = 15$), 3.9 (for $n = 25$) and 3.10 (for $n = 60$) display the empirical power curves of the alternative testing procedures, d , λ_0 , t_0 and τ_0 relating to the hypothesis $\phi = 0$. The curves show powers corresponding to a theoretical size $\alpha = .05$ for two-sided test; the one-sided power curves are given in the appendix.

The ordinate scale in Figures 3.8, 3.9 and 3.10 is linear in $\text{Arcsin}(\sqrt{\text{power}})$, going from 0° to 90° . In this scaling $\alpha = .05 \sim 12.9^\circ$, and the standard deviation of all points is approximately 1.4° , except when $\phi = 0$ where it is about 1° .

If the DM-test is credited with the power of d , then d still performs comparatively well in relation to Model A, in particular for $.0 < \phi < .9$. For negative and very large positive ϕ , the alternatives t_0 and τ_0 show greater power. In the context of Model B t_0 and τ_0 are superior to d . This relative deterioration demonstrates, that the mentioned shortcomings of the DM-test for models lacking a mean are substantive. It is noted that t_0 and τ_0 behave remarkably alike.

Throughout λ_0 is seemingly more powerful than t_0 and τ_0 , however this appears attributable to the fact that the size of the λ_0 -test is actually significantly larger than the sizes of t_0 and τ_0 , which as a whole do not differ significantly from $\alpha = .05$. In fact the λ_0 -curve

MODEL A, N=15

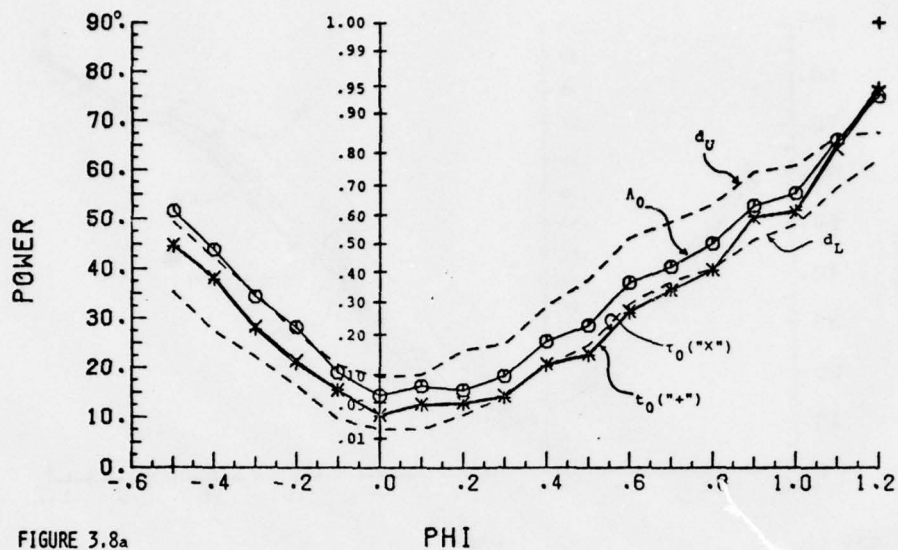


FIGURE 3.8a

PHI

113

MODEL B, N=15

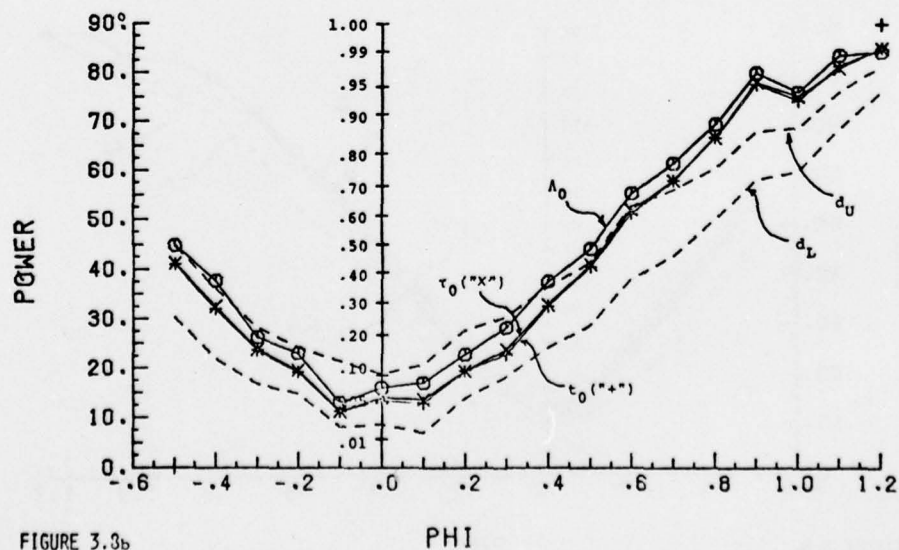


FIGURE 3.8b

PHI

114

MODEL A, N=25

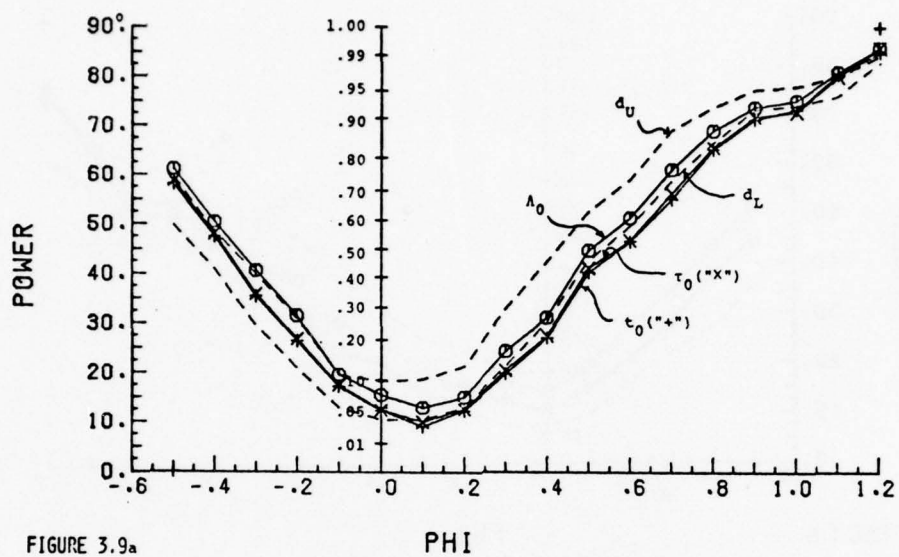


FIGURE 3.9a

15

MODEL B, N=25

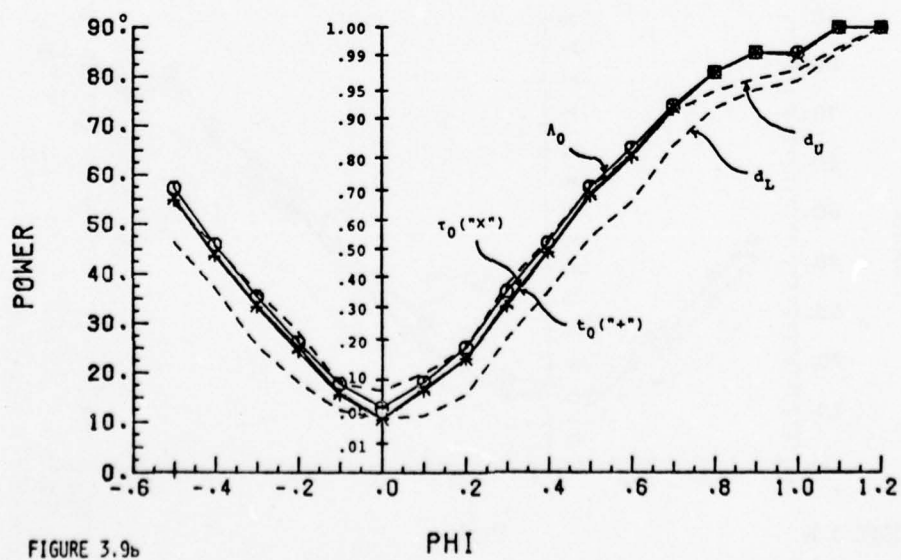
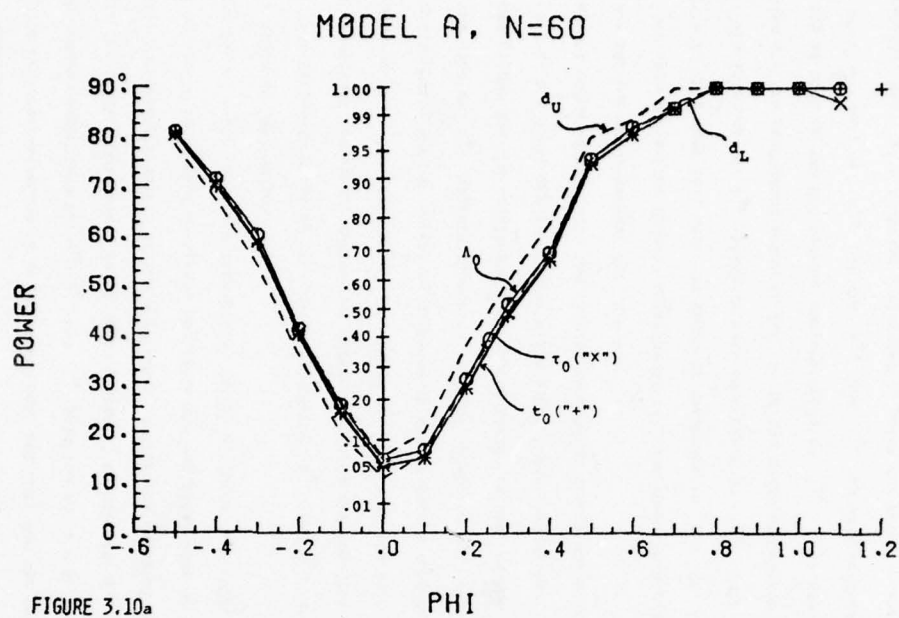
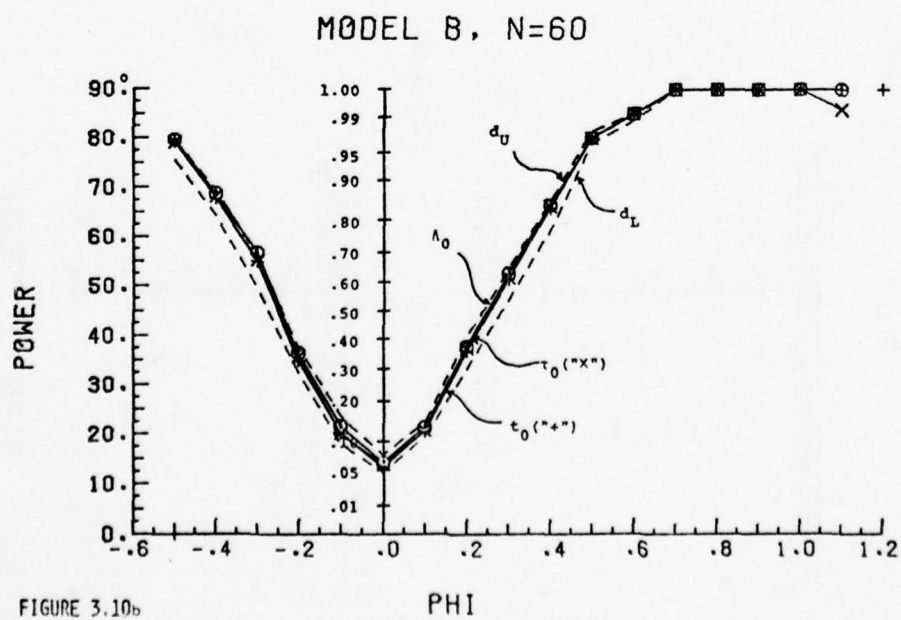


FIGURE 3.9b

16



117



118

looks much like those of t_0 and τ_0 , just shifted a bit upwards.

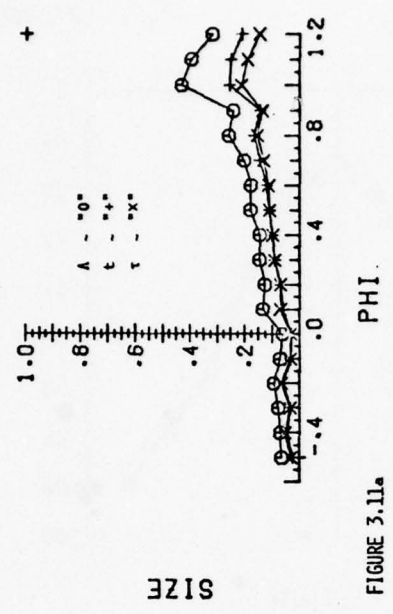
3.6.2 Size curves.

The size consideration is an important one, not only for the testing procedures symbolized by Λ_0 , t_0 and τ_0 relating to $\phi = 0$, but equally for their counterparts relating to a general hypothesis $\phi = \phi^*$. Actually the appropriateness of the approximate confidence interval defined by (3.47) is checked by seeing how close the empirical size of the χ^2 -test for $\phi = \phi^*$ is to its theoretical value α (here $\alpha = .05$), over a range of ϕ^* values.

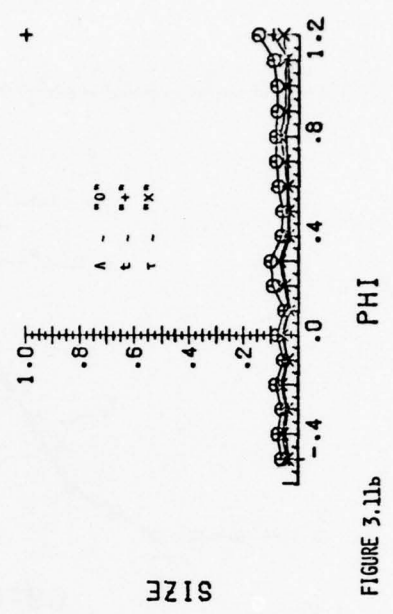
In relation to Model B, the empirical sizes of t_{ϕ^*} and τ_{ϕ^*} are almost identical throughout and do not differ significantly from the target value of $\alpha = .05$. For large ϕ^* , where the approximately non-informative prior for ϕ starts to distinguish itself from a uniform one, the size of τ_{ϕ^*} appears closer to $\alpha = .05$ than that of t_{ϕ^*} . The size graphs for the likelihood ratio based tests are also rather flat when $\phi^* < 1$, but their empirical levels lie significantly above the theoretical value of $\alpha = .05$, the displacement being most pronounced for $n = 15$, but still present for $n = 60$.

All size curves for Model A increase as ϕ^* increases to a peak at $\phi^* = 1$, then they fall off. In terms of closeness to $\alpha = .05$, Λ_{ϕ^*} is clearly inferior to t_{ϕ^*} , which in turn must yield to τ_{ϕ^*} . This explains why the confidence interval for ϕ in the examples Section 3.5 were not as wide as the HPD regions corresponding to τ_{ϕ^*} . The trouble when $\phi^* = 1$ not only for Λ_{ϕ^*} but for t_{ϕ^*} and τ_{ϕ^*} as well confirms empirically, that $\phi^* = 1$ being a singularity point for Model A cannot be handled like other points.

MODEL A, N=15



MODEL B, N=15



MODEL A, N=60

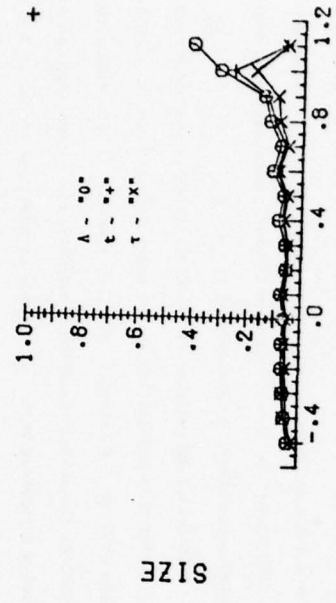


FIGURE 3.13a

MODEL B, N=60

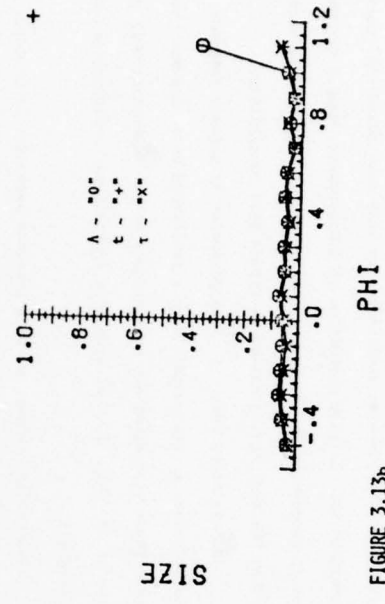


FIGURE 3.13b

MODEL A, N=25

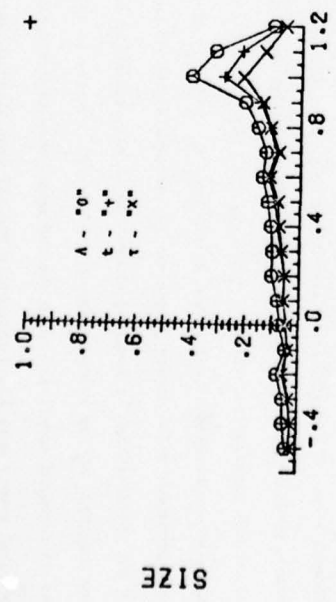


FIGURE 3.12a

MODEL B, N=25

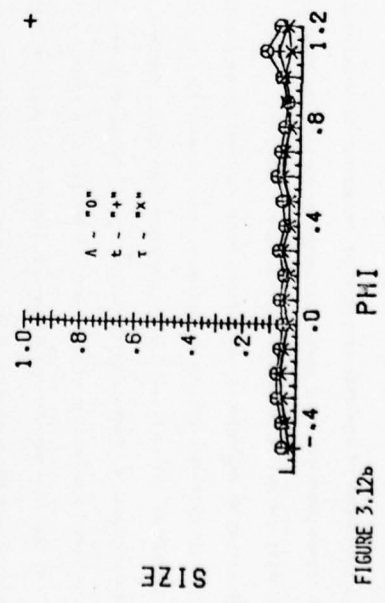


FIGURE 3.12b

Least squares regression is often applied in situations where it is feared that the observations are serially correlated. Since the usual techniques of analysis rest upon an assumption of independence it is a common practice under these circumstances to test for autocorrelation using the well known Durbin-Watson test. (This test is equivalent to obtaining an estimate $\hat{\phi}$ of a first order autoregressive parameter ϕ from the residuals left by an ordinary least squares fit; and then seeing whether $\hat{\phi}$ is significantly different from zero. Because the distribution of $\hat{\phi}$ depends on the matrix $X'X$ of independent variables, this is done by referring to bound for the desired significance point.)

In this Chapter some questions surrounding this practice have been approached from two points of view, namely using a likelihood approach and a Bayesian approach.

First the question may be raised whether the DW-test can be improved upon. A direct likelihood ratio approach yields an approximate significance test, Λ_0 , which, being asymptotic only, does not involve X . From a Bayesian point of view one can regard the inclusion or exclusion of some hypothesized value $\phi = \phi^*$ ($=0$) in, say, the $100(1-\alpha)\%$ HPD region as a statistical test. Test criterions t_0 and τ_0 of this kind were produced, which do not depend on X ; and while they do not require any justification by sampling properties it was studied by simulation how they and Λ_0 as alternatives to the DW-test, compare in terms of power in repeated sampling. For a studied model form involving a mean, the empirical powers were found comparable for all four alternatives; but the DW-test was somewhat less powerful in relation to a studied model form without a mean.

Secondly when serial dependence is to be expected, it may be questioned whether it actually makes sense to test a null hypothesis of independence. The DW-test seems often to have been abused to justify a convenient but unlikely assumption of independence, mainly for the purpose of proceeding as if, by a stroke of magic, independence were created by the test having failed to reject this possibility.

In situations where ϕ is a parameter of primary interest, inference about ϕ may be made from a likelihood or a Bayesian point of view. These approaches not only produce better (unconditional) estimates of ϕ than $\hat{\phi}$, but also give approximate confidence or HPD intervals for ϕ .

Ordinarily however ϕ is a nuisance parameter in relation to θ . As illustrated by the examples, proceeding in the fashion implied by the DW testing approach carries a penalty, as it may lead to drawing wrong inferences about the primary parameters θ , both in regard to location and to precision. This may happen when ϕ is different from zero but not detected significantly so, and inferences about θ subsequently are made conditional on $\phi = 0$. Or even if serial correlation is detected, but inferences about θ are made as if $\phi = \hat{\phi}$.

When serial correlation is expected, it would seem more sensible to allow for this possibility in the model and estimate ϕ simultaneously with θ . Minimizing $SS(\hat{\theta}, \hat{\phi})$, Equation (3.10), with respect to ϕ produces the maximum likelihood estimate, $\hat{\phi}$. From a sampling theory viewpoint inferences about θ may now be made conditional on $\phi = \hat{\phi}$ by substituting this estimate into the independence inducing transformation, c , and analyzing the transformed linear model using the usual techniques. From the Bayesian viewpoint this parallels approximation

of the marginal posterior distribution of $\hat{\phi}$ by the posterior distribution of $\hat{\phi}$ conditional on $\hat{\phi} = \hat{\phi}$.

Appendix. Additional results from the Monte Carlo study.

Figures 3A.1, 3A.2 and 3A.3 show the empirical onesided power curves of the alternative testing procedures, d , t_0 and τ_0 relating to the hypothesis $\phi = 0$, and at the theoretical level $\alpha = .05$.

Tables 3A.I, 3A.II and 3A.III list the averages and standard deviations, s , around the averages, of the estimators $\hat{\phi}$, $\hat{\phi}$ and $\hat{\phi}$ computed from the 400 replications for each model, A and B, for each sample size n and each ϕ -value. Overall $\hat{\phi}$ is clearly a much better estimator of ϕ than $\hat{\phi}$. In the Bayesian analysis the estimator $\hat{\phi}$ plays no prominent role except as a partial description of the marginal posterior distribution of ϕ . The difference between $\hat{\phi}$ and $\hat{\phi}$ is caused by the noninformative prior for ϕ (relative to a uniform one). It is noted, that on the average it does not make $\hat{\phi}$ overshoot ϕ .

MODEL A, N=15

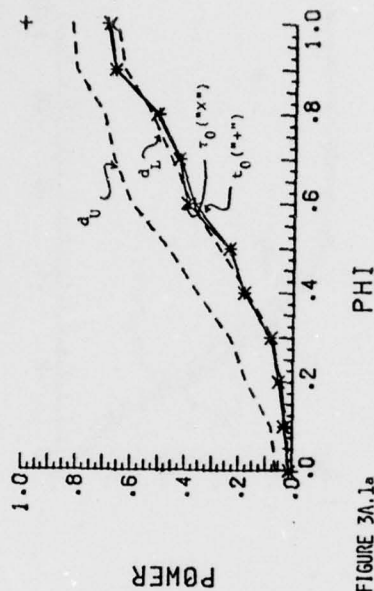


FIGURE 3A.1a

MODEL B, N=15

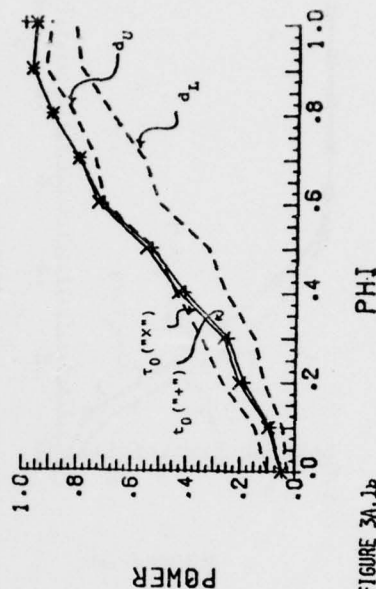
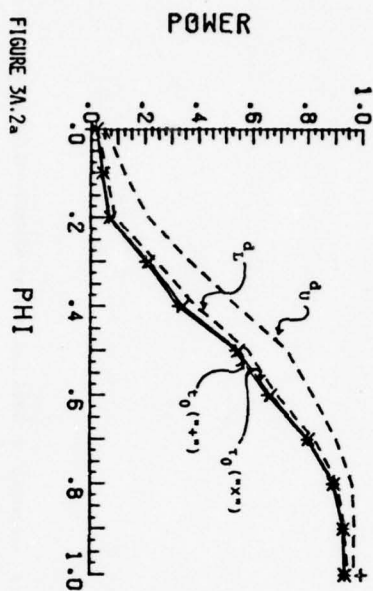
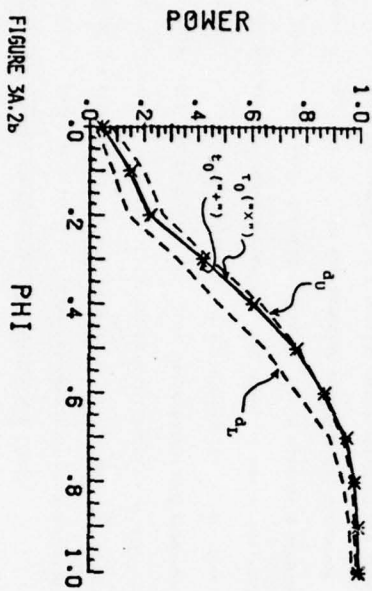


FIGURE 3A.1b

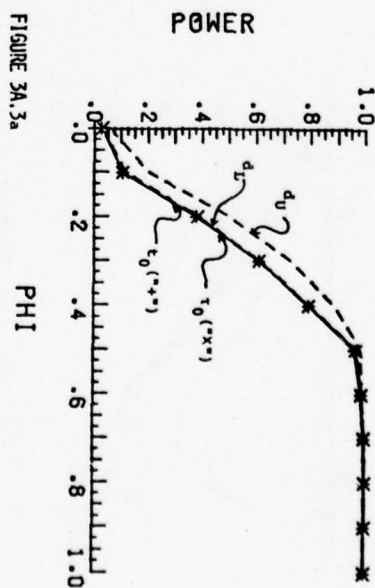
MODEL A, N=25



MODEL B, N=25



MODEL A, N=60



MODEL B, N=60

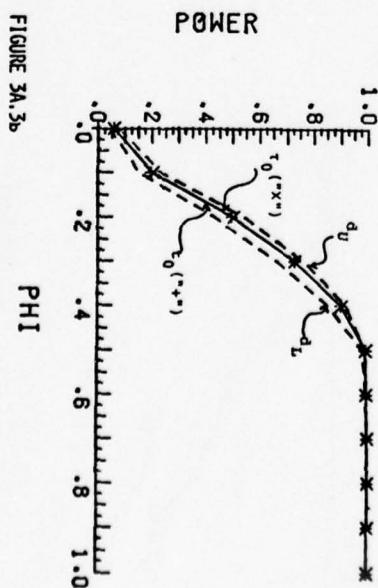


TABLE 3A.II

Model	ϕ	n	$\hat{\phi}$	s	$\hat{\phi}$	s	$\hat{\phi}$	s
A	.50	25	.438	.173	.504	.166	.528	.175
	.40	25	.360	.178	.419	.181	.434	.190
	.30	25	.272	.183	.330	.183	.346	.192
	.20	25	.182	.204	.233	.206	.244	.210
	.10	25	.108	.190	.157	.194	.165	.203
	.00	25	.045	.202	.083	.212	.086	.222
	.10	25	.044	.196	.088	.210	.090	.221
	.20	25	.126	.183	.176	.209	.184	.209
	.30	25	.201	.191	.264	.206	.278	.217
	.40	25	.283	.183	.359	.210	.374	.222
	.50	25	.367	.184	.437	.204	.460	.215
	.60	25	.434	.181	.523	.193	.554	.209
	.70	25	.510	.171	.612	.193	.655	.225
	.80	25	.584	.164	.712	.197	.768	.235
	.90	25	.654	.171	.825	.200	.880	.225
B	1.00	25	.689	.174	.967	.187	1.010	.211
	1.10	25	.784	.182	.987	.187	1.010	.211
	1.20	25	.840	.185	1.177	.097	1.187	.094
	.50	25	.403	.175	.474	.181	.494	.190
	.40	25	.333	.178	.385	.189	.403	.199
	.30	25	.242	.190	.283	.203	.296	.212
	.20	25	.147	.208	.178	.223	.187	.233
	.10	25	.076	.192	.095	.209	.100	.220
	.00	25	.010	.200	.010	.216	.011	.227
	.10	25	.080	.199	.094	.226	.094	.235
	.20	25	.159	.183	.186	.203	.195	.213
	.30	25	.247	.191	.292	.201	.305	.210
	.40	25	.323	.182	.381	.185	.399	.194
	.50	25	.405	.179	.474	.182	.486	.192
	.60	25	.475	.178	.556	.183	.582	.192
	.70	25	.544	.164	.661	.164	.692	.173
	.80	25	.634	.155	.744	.143	.780	.146
	.90	25	.718	.157	.856	.122	.896	.135
	1.00	25	.750	.152	.925	.123	.968	.135
	1.10	25	.807	.106	1.057	.093	1.085	.094
	1.20	25	.928	.068	1.188	.053	1.198	.054

TABLE 3A.I

Model	ϕ	n	$\hat{\phi}$	s	$\hat{\phi}$	s	$\hat{\phi}$	s
A	.50	15	.410	.216	.514	.216	.562	.262
	.40	15	.352	.219	.451	.232	.496	.269
	.30	15	.270	.221	.372	.249	.398	.281
	.20	15	.185	.235	.271	.283	.297	.289
	.10	15	.110	.227	.187	.253	.205	.280
	.00	15	.040	.225	.110	.254	.120	.277
	.10	15	.031	.240	.054	.241	.057	.310
	.20	15	.082	.255	.028	.245	.034	.337
	.30	15	.141	.267	.095	.270	.109	.325
	.40	15	.212	.253	.165	.305	.189	.342
	.50	15	.273	.251	.252	.366	.283	.356
	.60	15	.362	.244	.357	.391	.406	.365
	.70	15	.402	.256	.401	.368	.443	.344
	.80	15	.432	.251	.449	.395	.510	.424
	.90	15	.512	.233	.587	.293	.662	.359
B	1.00	15	.528	.243	.614	.304	.673	.344
	1.10	15	.628	.249	.624	.300	.695	.397
	1.20	15	.671	.234	1.045	.255	1.096	.272
	.50	15	.360	.225	.457	.243	.495	.267
	.40	15	.282	.220	.380	.261	.415	.290
	.30	15	.204	.225	.241	.271	.264	.298
	.20	15	.131	.236	.170	.290	.185	.317
	.10	15	.053	.217	.084	.248	.070	.293
	.00	15	.014	.215	.006	.281	.007	.308
	.10	15	.026	.234	.002	.273	.011	.328
	.20	15	.033	.224	.047	.261	.045	.306
	.30	15	.070	.249	.087	.276	.092	.294
	.40	15	.133	.269	.151	.282	.155	.284
	.50	15	.204	.251	.229	.241	.235	.265
	.60	15	.282	.251	.305	.241	.302	.265
	.70	15	.361	.227	.416	.205	.402	.273
	.80	15	.437	.221	.517	.206	.491	.251
	.90	15	.516	.201	.637	.163	.669	.176
	1.00	15	.637	.211	.702	.186	.769	.213
	1.10	15	.749	.207	1.026	.160	1.077	.168
	1.20	15	.824	.156	1.136	.154	1.179	.190

To Difference or not to Difference Data Series in
Linear Model Analysis: A Bayesian Study.

TABLE 3A.III

Model	ϕ	n	$\hat{\phi}$	s	$\hat{\phi}$	s	$\hat{\phi}$	s
A	-.50	60	-.476	.106	-.500	.107	-.510	.129
	-.40	60	-.382	.121	-.408	.120	-.416	.122
	-.30	60	-.285	.122	-.288	.122	-.294	.125
	-.20	60	-.206	.127	-.208	.126	-.212	.131
	-.10	60	-.104	.132	-.102	.134	-.104	.137
	.00	60	-.018	.121	-.035	.123	-.035	.125
	.10	60	.075	.132	.061	.133	.062	.136
	.20	60	.166	.119	.174	.123	.177	.125
	.30	60	.259	.121	.250	.124	.255	.127
	.40	60	.336	.124	.327	.127	.334	.130
	.50	60	.405	.117	.402	.123	.451	.126
	.60	60	.533	.116	.535	.121	.506	.125
B	-.50	60	-.624	.101	.631	.107	.606	.111
	-.40	60	.718	.099	.729	.105	.708	.110
	-.30	60	.790	.089	.818	.092	.841	.097
	-.20	60	.857	.090	.807	.088	.830	.115
	-.10	60	.842	.059	1.095	.029	1.098	.026
	.00	60	.462	.110	.489	.111	.408	.113
	.10	60	.366	.121	.391	.124	.309	.127
	.20	60	.239	.122	.310	.129	.316	.131
	.30	60	.191	.126	.202	.134	.205	.136
	.40	60	-.049	.132	-.092	.134	-.094	.137
	.50	60	-.001	.121	.001	.127	.001	.129
	.60	60	.091	.131	.099	.134	.101	.141
C	-.50	60	.199	.119	.214	.124	.218	.124
	-.40	60	.275	.123	.301	.125	.297	.127
	-.30	60	.354	.119	.378	.122	.375	.124
	-.20	60	.484	.115	.493	.117	.502	.119
	-.10	60	.553	.109	.565	.109	.597	.112
	.00	60	.644	.096	.683	.092	.699	.095
	.10	60	.737	.091	.761	.085	.800	.084
	.20	60	.819	.083	.875	.064	.896	.067
	.30	60	.901	.084	.972	.049	.995	.057
	.40	60	.972	.032	1.097	.022	1.099	.021
	.50	60	.972	.032	1.097	.022	1.099	.021
	.60	60	.972	.032	1.097	.022	1.099	.021

For the purpose of definiteness, consider some data analyzed by Coen, Comme & Kendall [1969] (CGK for short). In an attempt to forecast the Financial Times ordinary share index, they relate this dependent variable, y_i , detrended, to two detrended "leading" independent variables $x_{1,i}$ and $x_{2,i}$, viz. U.K. car production and the Financial Times commodity index (their equation number 7):

$$y_i = \theta_1 + \theta_2 x_{1,i} + \theta_3 x_{2,i} + \theta_4 i + N_i \quad (4.1)$$

where $x_{1,i}$ is the observed value 6 time periods (quarters) previous to y_i , and $x_{2,i}$ has a lead time of 7 periods.

This data set raises a number of questions. First of all whether the noise N_i might be serially correlated and not independent as the CGK analysis tacitly assumed. Secondly how this serial correlation, if indeed present, affects the inferences about the regression parameters θ_j , where CGK's estimates of θ_2 and θ_3 were found highly significantly different from zero.

Both these questions were addressed by Box and Newbold [1971] (BN), who demonstrated, that the noise was in fact highly serially dependent, and that taking the autocorrelation into account the claimed dependencies of y on x_1 and x_2 disappeared. BN entertained several alternative noise structures, the most prominent of which were a first order autoregressive (AR-1) model, and a first order integrated moving average (IMA-1) model (i.e. a noise model for which the first difference is a first order moving average process). This latter model fits the data only slightly less well than the AR-1 model, but BN did not

commit themselves to any preferred form over the other, since it made little difference in demonstrating their point. This does however leave a third question, namely whether differencing should be carried out or not.

Writing the linear model

$$\begin{matrix} \bar{y} \\ n \times 1 \end{matrix} = \begin{matrix} X \\ n \times p \end{matrix} \begin{matrix} \bar{\theta} \\ p \times 1 \end{matrix} + \begin{matrix} \bar{N} \\ n \times 1 \end{matrix} \quad (4.2)$$

the noise \bar{N} may in great generality be considered generated by an autoregressive integrated moving average (ARIMA) process of order (p, d, q) , see Box and Jenkins [1970]. The presence of the deterministic linear model parameters $\bar{\theta}$ does not impede a time series modelling of the noise, as the linear and the time series parameters may be estimated simultaneously, and the latter are (asymptotically) distributed as if the deterministic model part did not exist. In such a general approach, the question of differencing takes the form of deciding the order d in the ARIMA model.

The scope of the present investigation is limited to situations, where it is impossible (due to a relatively small number of observations), impractical or perhaps unnecessary to employ more than one first order autoregressive parameter to allow for autocorrelation.

How such a model may be analyzed was the topic of Chapter 2, specifically it was studied how Bayesian inference is made about the model parameters $\bar{\theta}$ and ϕ_0 when the noise e_i of the observations (original or perhaps differenced) follows the AR-1 scheme

$$\begin{cases} e_1 = M_0 + \epsilon_1 \\ e_i = \phi_0 e_{i-1} + \epsilon_i \end{cases} \quad i = 2, 3, \dots, n \quad (4.3)$$

where ϕ_0 is the autoregressive parameter, M_0 is a starting parameter included to accommodate explosive, ($1 < \phi_0$) as well as stationary situations and the ϵ_i 's are i.i.d. $N(0, \sigma^2)$, (also independent of the p series of predictor variables \bar{X}).

Two ways of assessing whether or not differencing appears justified in light of the data will now be explored from a Bayesian point of view.

Within the AR-1 noise structure, differencing corresponds to a particular choice of the autoregressive parameter ϕ_0 , namely $\phi_0 = 1$, so it may be wondered whether the data support this value (Section 4.2).

Since the question of differencing is intimately related to the existence (or lack) of a fixed mean for the noise, it may usefully be regarded as one of determining whether one degree of freedom is better *expended on differencing the series* (effectively losing one observation) rather than on including a mean, while in either case serial correlation is being allowed for by one parameter. In other words, the testimony of the data concerning the appropriateness of differencing here takes the form of relative support of two alternative models with a common deterministic part but with different noise models. (Section 4.3).

The CGK data are reanalyzed in detail (Section 4.4) in particular the results concerning whether to difference or not are illustrated.

4.1 Joint posterior distribution of $\{\theta, \phi_0\}$

Analyzing data sets like the CGK data from a Bayesian point of view makes it very clear, why misleading conclusions may be reached about the linear parameters $\bar{\theta}$ when serial correlation is ignored.

Assuming that the data \bar{y} may adequately be represented by the linear model

$$\tilde{y} = \frac{1}{n} \tilde{\theta}_1 + \frac{1}{n} \tilde{\theta}_2 (1) + \tilde{\theta} \quad (4.4)$$

where

$$\tilde{X} = \begin{bmatrix} 1 & \tilde{X} \\ \tilde{X} & \tilde{X} \end{bmatrix}; \quad \tilde{\theta} = \begin{bmatrix} \theta \\ \theta \end{bmatrix} \quad (4.5)$$

with $\text{rank}(\tilde{X}) = p$, and where $\tilde{\theta}$ follows the AR-1 scheme (4.3); then it was seen in Chapter 2, that the density function for \tilde{y} is:

$$\begin{aligned} p(\tilde{y} | M_0, \tilde{\theta}, \phi_0, \sigma) &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2}\sigma^2 (\tilde{y} - \tilde{X}\tilde{\theta})' C (\tilde{y} - \tilde{X}\tilde{\theta}) - (\tilde{y}_1 - \tilde{X}_1 \tilde{\theta}_1)^2\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2}\sigma^2 (SS_0(\tilde{\theta}, \phi_0) + (e_1 - M_0)^2)\right) \quad (4.6) \end{aligned}$$

where

$$SS_0 = C_0' C_0; \quad C_0 = \begin{bmatrix} -\phi_0 & 1 & 0 \\ 0 & 0 & \phi_0 \end{bmatrix} \quad (4.7)$$

so that

$$SS_0(\tilde{\theta}, \phi) = (C_0 \tilde{y} - S_0 \tilde{X} \tilde{\theta})' (C_0 \tilde{y} - S_0 \tilde{X} \tilde{\theta}) \quad (4.8)$$

(and where \tilde{X}_1 is the first row of \tilde{X}). The subscript "0" is used to indicate, that the AR-1 noise relates to the "not differenced" data.

Of course the expression (4.6) is the likelihood function for the parameters $\tilde{\theta} = \{M_0, \tilde{\theta}, \phi_0, \sigma\}$ of the model (4.4) when \tilde{y} is considered known. It is noted, that maximizing the likelihood function is equivalent to minimizing the (residual) sum of squares of the estimated $\tilde{\theta}$'s, $SS_0(\tilde{\theta}, \phi_0)$; and the minimum is $SS_0(\hat{\theta}, \hat{\phi}_0)$.

It was argued in Chapter 2, that an approximately noninformative prior distribution for $\tilde{\theta}$ is of the form:

$$p(M_0, \tilde{\theta}, \phi_0, \sigma) = |X' C X|^{1/2} \sigma^{-1} p(\phi_0) \quad (4.9)$$

where $p(\phi_0)$ is a locally noninformative prior for ϕ_0 . The argument leading to (4.9) applies only when $X' C X$, or equivalently $C X$, has

full column rank p . Excluding the distinct singularity points $\phi = \hat{\phi}$, if any, where $\text{rank}(C X) < p$ (in fact equal to $p-1$), the joint posterior for $\{M_0, \tilde{\theta}, \phi_0, \sigma\}$, found by multiplying (4.6) by (4.9), was integrated over $\{M_0, \tilde{\theta}, \sigma\}$ to yield the marginal posterior distribution of ϕ_0 :

$$p(\phi_0 | \tilde{y}) = p(\phi_0) (SS_0(\hat{\theta}, \phi_0))^{-\frac{n-p-1}{2}} \quad (4.10)$$

where

$$\hat{\theta} = \hat{\theta}(\phi_0) = (X' C X)^{-1} X' C \tilde{y} \quad (4.11)$$

Omitting the distinct point(s) $\phi_0 = \hat{\phi}$ (if any), (4.10) is normalized as if it was a continuous distribution. Specifically it was suggested in Chapter 2 how $p(\phi_0 | \tilde{y})$ may be approximated by a t-distribution.

The exclusion of possible singularity points is only a matter of formality as far as marginal inference about $\tilde{\theta}$ is concerned.

Conditional on ϕ_0 ($\neq \hat{\phi}$) the distribution of $\tilde{\theta}$ a posteriori is $p(\tilde{\theta} | \phi_0, \tilde{y}) \sim t(\hat{\theta}, S^2 (X' C X)^{-1}, v)$ (4.12)

with

$$v S^2 = SS_0(\hat{\theta}, \phi_0); \quad v = n-p-1 \quad (4.13)$$

in particular

$$p(\tilde{\theta} | \phi_0, \tilde{y}) \sim t(\hat{\theta}_j, S^2 D_{jj}, v) \quad (4.14)$$

where D_{jj} comes from

$$[X' C X]^{-1} = \begin{bmatrix} D_{11} & & \\ & \ddots & \\ & & D_{jj} & \dots & \\ & & & \ddots & \\ & & & & D_{pp} \end{bmatrix} \quad \text{sym} \quad (4.15)$$

Now the joint posterior distribution of $\{\theta_j, \phi_0\}$ may conveniently be computed as

$$p(\theta_j, \phi_0 | \tilde{y}) = p(\theta_j | \phi_0, \tilde{y}) p(\phi_0 | \tilde{y}) \quad (4.16)$$

and plotting contours of this joint distribution, $j = 1, 2, \dots, p$, allows a clear appreciation of the dependence of θ_j on ϕ_0 . The levels of the contours corresponding to specified interior probability contents may be determined numerically. As an approximation it may be used (see for example Box and Tiao [1973] p. 94) that for large n $p(\theta_j, \phi_0 | y)$ tends towards normality, so that

$$-2 \ln \frac{P_0(\theta_j, \phi_0)}{P_0(\hat{\theta}_j, \hat{\phi}_0)} \sim \chi_2^2 \quad (4.17)$$

which (in the 2 dimensional case) may be expressed as:

$$\frac{P_0(\theta_j, \phi_0 | y)}{P_0(\hat{\theta}_j, \hat{\phi}_0 | y)} = e^{-\frac{1}{2} \chi_2^2(\alpha)} = \alpha \quad (4.18)$$

where $\hat{\phi}_0$ is determined by the identity

$$SS_0(\hat{\theta}_j, \hat{\phi}_0) = \min_{\phi_0} SS(\hat{\theta}_j, \phi_0) \quad (4.19)$$

and where $(\hat{\theta}_j^0, \hat{\phi}_0^0)$ is a point on the boundary of the 100(1- α)% highest posterior density (HPD) region.

4.2 Posterior density of the AR-1 parameter at unity.

For models like (4.4) involving a mean, the point $\phi_0 = 1$ is a singularity point. Thus the plausibility of ϕ_0 being unity (i.e. whether or not a differencing appears supported by the data) cannot be assessed by looking at $p(\phi_0 | y)$ (4.10) as a continuous distribution, since the prior (4.9) does not find justification at this point, and the integration leading to (4.10) does not go through.

The cause of the singularity at $\phi_0 = 1$ is the vanishing of the mean θ , at that point. Specifically (4.4) degenerates to

$$y = \sum_{i=0}^{\infty} \theta_i(1) \tilde{z}(1) + e(\phi_0 = 1) \quad (4.20)$$

for $\phi_0 = 1$; as namely $e_1(\phi_0 = 1)$, and therefore y_1 , does not possess a mean. Hence conditional on $\phi_0 = 1$ we have

$$\hat{\theta}_{(1)} = (\sum_{i=1}^n D X_{(1)}^i)^{-1} X_{(1)}' D Y \quad (4.21)$$

where

$$D = \begin{bmatrix} d_1' & d_2' & \dots & d_n' \\ 0 & 0 & \dots & 0 \end{bmatrix} \quad (4.22)$$

and

$$P_0(\theta_{(1)} | \phi_0 = 1, Y) \sim t_{p-1}(\hat{\theta}_{(1)}, s^2(X_{(1)}' D X_{(1)})^{-1}, v) \quad (4.23)$$

with

$$v s^2 = SS_0(\hat{\theta}_{(1)}, \phi_0 = 1) = (dy - d \sum_{i=1}^n X_{(1)}^i)' (dy - d \sum_{i=1}^n X_{(1)}^i)' : v = n-p \quad (4.24)$$

It is noted that $SS_0(\hat{\theta}_{(1)}, \phi_0 = 1)$ is in general not equal to $\lim_{\phi_0 \rightarrow 1} SS_0(\hat{\theta}, \phi_0)$. It was seen in Chapter 3 that the latter quantity, which we shall refer to as $SS_0(\hat{\theta}, \phi_0 = 1)$, may be computed as the residual sum of squares resulting from replacing, when $\phi_0 = 1$, the redundant parameter θ_1 by a new parameter θ_1^* serving as a mean for the differenced series, i.e. implying that the model incorporates a deterministic linear trend. So when this artificial trend is not incorporated, then the residual sum of squares, SS_0 , corresponding to the model (4.4) makes a discontinuous upwards jump at $\phi_0 = 1$ as a result of the model's loss of one parameter at that point. Thus if the prior for ϕ_0 is taken to be continuous through $\phi_0 = 1$ then the posterior density function for ϕ_0 will not in general be continuous at that point.

In order to determine a prior distribution at $\phi_0 = 1$ for $\{M_0, \hat{\theta}_{(1)}, \phi_0, \sigma\}$ which is known up to the same multiplicative constant as $p(M_0, \hat{\theta}, \phi_0, \sigma)$ of (4.9), we shall argue the following way: Conditional on σ and $\hat{\theta}_0$ the prior for the linear parameters $\hat{\theta}_{(1)}$ and M is taken locally uniform proportional to k . To determine k we suppose that σ is known, and look at the posterior densities of ϕ_0 at $\phi_0 = 1$ and $\phi_0 = \phi_0' (\neq 1)$. We shall then consider the prior for

$(M_{0,0}^{\theta} | \phi_0 = 1)$ being consistent with $p(M_{0,0}^{\theta} | \phi_0 = \phi_0')$ in the sense

$$\lim_{\phi_0' \rightarrow 1} E \ln p(\phi_0 = \phi_0' | y) = E \ln p(\phi_0 = 1 | y). \quad (4.25)$$

We find

$$p(\phi_0 = \phi_0' | y) = \int_{M_{0,0}^{\theta}} p(M_{0,0}^{\theta} | \phi_0 = \phi_0') p(\phi_0 = \phi_0' | y) p(M_{0,0}^{\theta} | \phi_0 = 1) dM_{0,0}^{\theta} \\ = p(\phi_0 = \phi_0') (2\pi\sigma^2)^{-\frac{n-p-1}{2}} \exp(-\frac{1}{2} \sigma^{-2} SS_0(\hat{\theta}, \phi_0 = \phi_0')) \quad (4.26)$$

and

$$p(\phi_0 = 1 | y) = \int_{M_{0,0}^{\theta}} p(M_{0,0}^{\theta} | \phi_0 = 1) p(\phi_0 = 1 | y) p(M_{0,0}^{\theta} | \phi_0 = 1) dM_{0,0}^{\theta} \\ = p(\phi_0 = 1) |X_{\infty}^{\theta}| \frac{D X_{\infty}^{\theta}}{X_{\infty}^{\theta}} |^{-1/2} k(2\pi\sigma^2)^{-\frac{n-p}{2}} \exp(-\frac{1}{2} \sigma^{-2} SS_0(\hat{\theta}, \phi_0 = 1)) \quad (4.27)$$

where it is understood in (4.26) and (4.27), that the proportionality constants are the same. Taking expectation of the logarithm of these expressions and equating gives

$$\ln p(\phi_0 = \phi_0') - \frac{n-p-1}{2} \ln(2\pi\sigma^2) - \frac{1}{2} (n-p-1) \\ = \ln p(\phi_0 = 1) \ln(|X_{\infty}^{\theta}| \frac{D X_{\infty}^{\theta}}{X_{\infty}^{\theta}} |^{-1/2} k) - \frac{n-p}{2} \ln(2\pi\sigma^2) - \frac{1}{2} (n-p) \quad (4.28)$$

Letting $\phi_0' \rightarrow 1$ the prior densities at $\phi_0 = \phi_0'$ and $\phi_0 = 1$ cancel, leaving

$$k \cdot |X_{\infty}^{\theta}| \frac{D X_{\infty}^{\theta}}{X_{\infty}^{\theta}} |^{-1/2} e^{1/2} (2\pi\sigma^2)^{1/2} \quad (4.29)$$

hence we are led to choose

$$p(M_{0,0}^{\theta} | \phi_0 = \phi_0', \sigma) = p(M_{0,0}^{\theta} | \sigma, \phi_0 = 1) p(\sigma) p(\phi_0) |_{\phi_0 = 1} \\ = |X_{\infty}^{\theta}| \frac{D X_{\infty}^{\theta}}{X_{\infty}^{\theta}} |^{-1/2} e^{1/2} (2\pi\sigma^2)^{1/2} \sigma^{-1} p(\phi_0 = 1). \quad (4.30)$$

The posterior density $p(\phi_0 | y)$ at $\phi_0 = 1$ is now determined from considering

$$\frac{p(\phi_0 = 1 | y)}{p(\phi_0 = \phi_0' | y)} = \frac{p(\phi_0 = 1)}{p(\phi_0 = \phi_0')} \frac{\int_{M_{0,0}^{\theta}} I_{\phi_0 = 1} (M_{0,0}^{\theta} | \phi_0 = 1, \sigma) dM_{0,0}^{\theta}}{\int_{M_{0,0}^{\theta}} I_{\phi_0 = \phi_0'} (M_{0,0}^{\theta} | \phi_0 = \phi_0', \sigma) dM_{0,0}^{\theta}} \\ = \frac{p(\phi_0 = 1)}{p(\phi_0 = \phi_0')} \frac{e^{1/2} (SS_0(\hat{\theta}, \phi_0 = 1))^{-\frac{n-p-1}{2}}}{(SS_0(\hat{\theta}, \phi_0 = \phi_0'))^{-\frac{n-p-1}{2}}} \quad (4.31)$$

with

$$I_{\phi_0 = 1} (M_{0,0}^{\theta} | \phi_0 = 1, \sigma) =$$

$$|X_{\infty}^{\theta}| \frac{D X_{\infty}^{\theta}}{X_{\infty}^{\theta}} |^{-1/2} e^{1/2} (2\pi\sigma^2)^{-\frac{n-1}{2}} \sigma^{-1} \exp(-\frac{1}{2} \sigma^{-2} (SS_0(\hat{\theta}, \phi_0 = 1) + (e_1 - M_1)^2))$$

and

$$I_{\phi_0 = \phi_0'} (M_{0,0}^{\theta} | \phi_0 = \phi_0', \sigma) =$$

$$|X_{\infty}^{\theta}| \frac{D X_{\infty}^{\theta}}{X_{\infty}^{\theta}} |^{-1/2} e^{1/2} (2\pi\sigma^2)^{-\frac{n-1}{2}} \sigma^{-1} \exp(-\frac{1}{2} \sigma^{-2} (SS_0(\hat{\theta}, \phi_0 = \phi_0') + (e_1 - M_1)^2)).$$

Now letting $\phi_0' \rightarrow 1$, we find the following expression for the density at $\phi_0 = 1$:

$$p(\phi_0 = 1 | y) = \lim_{\phi_0' \rightarrow 1} p(\phi_0 | y) e^{1/2} \left(\frac{SS_0(\hat{\theta}, \phi_0 = 1)}{SS_0(\hat{\theta}, \phi_0 = \phi_0')} \right)^{-\frac{n-p-1}{2}} \quad (4.32)$$

The discrete jump of $p(\phi_0 | y)$ at $\phi_0 = 1$ may be upwards or downwards depending on the last factor in (4.32); since it is less than or equal to 1, the largest possible upwards jump is by a factor $e^{1/2}$; 1.65.

Other singularity point $\phi_0 = \phi_0'$ (if any), may obviously be taken care of similarly.

4.3 Expanding one degree of freedom on differencing rather than on a fixed mean?

In this section the question of whether some autocorrelated series should perhaps be differenced prior to a linear model analysis is approached from the point of view of determining to what extent the data are more consistent with a model having a fixed mean and an AR-1 noise (4.3), or with a model which includes no mean and has a nonstationary noise component. In the Bayesian framework the testimony of the data concerning this question may be extracted in terms of posterior probabilities for two alternative models with a common deterministic part, but with alternative noise models, both allowing for serial correlation by one parameter.

In this formulation the problem is similar to the one known as model discrimination, Box & Hill [1967], Box & Henson [1969], Kanemasu [1973]. Their problem was discrimination among alternative deterministic models with identical noise components; in the present problem the discrimination is between alternative noise models for the same deterministic model.

Singularity points require no special attention in this analysis, in particular the model (4.4) involves, with probability one, a fixed mean θ_1 .

4.3.1 The alternative likelihood function

A natural alternative, \tilde{z} , to the AR-1 noise structure z (4.3), consists in letting the noise of the differenced series follow an AR-1 process, i.e.

$$\tilde{d} \tilde{y} = d \tilde{x} \tilde{z} \tilde{z}(1) + d \tilde{z} \tilde{z} \quad (4.33)$$

or

$$\tilde{y}(1) = \tilde{w} \tilde{z}(1) + \tilde{z}(1) \quad (4.34)$$

where

$$\tilde{y}(1) = \begin{bmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_n \end{bmatrix} = d \tilde{y} ; \tilde{w} = \begin{bmatrix} \tilde{w}_1 \\ \vdots \\ \tilde{w}_n \end{bmatrix} = d \tilde{z}(1) \quad (4.35)$$

and

$$\begin{cases} \tilde{e}_2 = \tilde{m}_1 + \tilde{e}_2 \\ \tilde{e}_i = \phi_1 \tilde{e}_{i-1} + \tilde{e}_i \quad i = 3, 4, \dots, n \end{cases} \quad (4.36)$$

where the \tilde{e}_i 's are i.i.d. $N(0, \sigma^2)$.

Hence the noise vector \tilde{z} relating to the original observations can be represented by the stochastic model

$$(1 - \phi_1 B) \nabla \tilde{z} = \tilde{e} \quad (4.37)$$

where $\nabla = (1 - B)$. (4.37) is an ARIMA model of order (1,1,0), which incidentally in the present formulation also covers explosive situations, $1 < \phi_1$.

Introducing H_0 and H_1 for the two alternative models:

$$H_0 : \tilde{y} = X \tilde{\theta} + \tilde{e} \quad (4.38)$$

$$H_1 : \tilde{y} = X \tilde{z}(1) + \tilde{f} \quad (4.39)$$

the distribution function of \tilde{y} in relation to H_0 is given by Equation (4.6), since of course (4.38) is a repetition of (4.4).

In order to derive the distribution function of \tilde{y} relative to H_1 , we shall set $\tilde{v}_1 = \tilde{y}_1$ and define the parameter $\tilde{\theta}_0$ so that

$$\tilde{v} = \begin{bmatrix} \tilde{v}_1 \\ \tilde{v}(1) \end{bmatrix} = \begin{bmatrix} \tilde{\theta}_0 \\ \tilde{w} \tilde{z}(1) \end{bmatrix} + \begin{bmatrix} \tilde{e}_1 \\ \tilde{z}(1) \end{bmatrix} \quad (4.40)$$

where $\hat{\epsilon}_1$ is distributed jointly with the $\hat{\epsilon}_i$'s of (4.36) as

$$p(\hat{\epsilon}_1) = (2\pi\sigma^2)^{-n/2} \exp(-\frac{1}{2}\sigma^{-2} \sum_{i=1}^n \hat{\epsilon}_i^2). \quad (4.41)$$

From (4.41) and (4.36) we find

$$p(\hat{\epsilon}_1, \hat{\epsilon}_2) = (2\pi\sigma^2)^{-n/2} \exp(-\frac{1}{2}\sigma^{-2} (\hat{\epsilon}_1^2 + (\hat{\epsilon}_2 - M_1)^2 + \hat{\epsilon}_1^2)) \quad (4.42)$$

where

$$C_{n1} = C_{n1}^{\epsilon_1, \epsilon_2} = \begin{bmatrix} -\hat{\epsilon}_1 & 1 & 0 \\ 0 & 0 & -\hat{\epsilon}_1 \\ 0 & 0 & 1 \end{bmatrix} \quad (4.43)$$

Now (4.42) in conjunction with (4.40) yields

$$p(y|M_1, \theta, \phi_1, \sigma) = \frac{1}{(n-2)!} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(-\frac{1}{2}\sigma^{-2} ((v_1 - M_1)^2 + (v_2 - M_1)^2 + (v_1 - \theta)^2)) \quad (4.44)$$

and since the transformation

$$y = \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \hat{\epsilon}_1 \\ \hat{\epsilon}_2 \\ \hat{\epsilon}_3 \end{bmatrix} \quad (4.45)$$

has unit Jacobian, the right hand side of (4.44) is also the density

function of y :

$$p(y|M_1, \theta, \phi_1, \sigma) = (2\pi\sigma^2)^{-n/2} \exp(-\frac{1}{2}\sigma^{-2} (SS_1(\theta, \phi_1) + (v_1 - M_1)^2 + (v_1 - \theta)^2)) \quad (4.46)$$

where

$$SS_1(\theta, \phi_1) = (v_1 - M_1)^2 + C_1(v_1 - M_1) \quad (4.47)$$

much like (4.9). When y is considered known, (4.46) is of course the likelihood function for the parameters $(M_1, \theta, \phi_1, \sigma)$, and it is maximized by minimizing the (residual) sum of squares of the estimated $\hat{\epsilon}_i$'s, $SS_1(\theta, \phi_1)$; the minimum is $SS_1(\hat{\theta}, \hat{\phi}_1)$.

Since H_0 and H_1 have the same number of parameters, $p+3$, we may formally equate

$$\begin{cases} p(y|\hat{\epsilon}_1, H_0) = p(y|M_0, \theta, \phi_0, \sigma) \\ p(y|\hat{\epsilon}_1, H_1) = p(y|M_1, \theta, \phi_1, \sigma) \end{cases} \quad (4.48)$$

and these distributions are known exactly including their normalizing constants.

4.3.2 The complementing priors.

Before posterior model probabilities can be computed, prior distributions for y complementing the likelihood functions must be produced.

As far as $\hat{\epsilon}_0 = p(y|\hat{\epsilon}_0, H_0)$ is concerned, we have (4.10):

$$p(\hat{\epsilon}_0|H_0) = p(M_0, \theta, \phi_0, \sigma) = |X' C_0 X|^{1/2} \sigma^{-1} p(\phi_0) \quad (4.49)$$

An approximately noninformative prior complementing the nonstationary alternative likelihood function $\hat{\epsilon}_1 = p(y|\hat{\epsilon}_1, H_1)$ is derived as follows. On an assumption of independence a priori between the linear parameter θ measuring an initial level of the series, and the remaining parameters, we may write

$$p(\hat{\epsilon}_1|H_1) = p(M_1, \theta, \phi_1, \sigma) = p(M_1, \theta, \phi_1, \sigma) p(\theta) \quad (4.50)$$

where now the prior $p(M_1, \theta, \phi_1, \sigma)$ may be viewed as complementing the conditional likelihood

$$\begin{aligned}
& I_1(M_1, \theta_1, \phi_1, \sigma_1^2 | \theta_1, \phi_1, \sigma_1^2) = I_1(M_1, \theta_1, \phi_1, \sigma_1^2 | \theta_1, \phi_1, \sigma_1^2) \\
& = p(\theta_1 | \theta_1, \phi_1, \sigma_1^2) p(\phi_1 | \theta_1, \phi_1, \sigma_1^2) p(\sigma_1^2 | \theta_1, \phi_1, \sigma_1^2) \\
& = (2\pi\sigma_1^2)^{-\frac{n-1}{2}} \exp\left(-\frac{1}{2\sigma_1^2} \left((v_1 - \theta_1)^2 + (v_2 - \theta_1)^2 + \dots + (v_n - \theta_1)^2 \right) \right) \\
& \quad \cdot \exp\left(-\frac{1}{2\sigma_1^2} \left((v_1 - \theta_1)^2 + (v_2 - \theta_1)^2 + \dots + (v_n - \theta_1)^2 \right) \right) \quad (4.51)
\end{aligned}$$

which has the familiar AR-1 form as before, (4.6). Hence on the same arguments (presented in Chapter 2) that led to the adoption of the approximately noninformative prior (4.49) in conjunction with H_0 , and assigning a locally uniform prior for the location parameter θ_0 , we find

$$p(\psi | H_1) = \left| \frac{W' C_1 W}{n} \right|^{1/2} \sigma^{-1} p(\phi_1) \quad (4.52)$$

where as before $p(\phi_1)$ is a locally approximately noninformative prior for ϕ_1 .

Equations (4.49) and (4.52) specify the priors $p(\psi | H_m)$, $m = 0, 1$, up to a multiplicative constant only. In order to compute posterior model probabilities, it is necessary to know the ratio of these proportionality constants; in fact prior model probabilities need also be stated. First utilizing the (usual) assumption of prior independence between σ and the remaining parameters we may factorize:

$$p(\psi | H_0) = p(M_0, \theta_0, \phi_0, \sigma) = p(\theta_0 | M_0, \phi_0) p(M_0 | \phi_0) p(\phi_0) p(\sigma) \quad (4.53)$$

and

$$p(\psi | H_1) = p(M_1, \theta_1, \phi_1, \sigma) = p(\theta_1 | M_1, \phi_1) p(M_1 | \phi_1) p(\phi_1) p(\sigma). \quad (4.54)$$

Now further progress can be made, if we make the following assumptions:

1) The constants of proportionality in

$$p(\theta | M_0, \phi_0) = \left| \frac{W' C_0 W}{n} \right|^{1/2} \quad (4.55)$$

and

$$p(\theta | M_1, \phi_1) = \left| \frac{W' C_1 W}{n} \right|^{1/2} \quad (4.56)$$

are the same.

It is noted, that conditionally on M_m, ϕ_m , $m = 0, 1$, the two alternative models are equivalent to

$$H_0 | M_0, \phi_0 : z_0 = \theta_0 + \epsilon \quad (4.57)$$

with $z_0 = \theta_0 + \epsilon$; $\epsilon_0 = \epsilon_0 + \epsilon$

$$H_1 | M_1, \phi_1 : z_1 = \theta_1 + \epsilon \quad (4.58)$$

$$\text{with } z_1 = \begin{bmatrix} v_1 \\ c_1 v_1 \end{bmatrix}; \quad \theta_1 = \begin{bmatrix} 1 & 0' \\ 0 & C_1 W \end{bmatrix}$$

and models (4.57) and (4.58) are ordinary linear models of p linear parameters, and i.i.d. $N(0, \sigma^2)$ noise.

Now if we employ to this situation the argument of Kanemasu [1973], that the increase in information (measured by entropy) about the parameters θ and (θ_1, ϕ_1) from n observations, is expected to be the same a priori; then his result says, that indeed

$$\frac{p(\theta | M_0, \phi_0)}{p(\theta_1 | M_1, \phi_1)} = \frac{\left| \frac{W' C_0 W}{n} \right|^{1/2}}{\left| \frac{W' C_1 W}{n} \right|^{1/2}} = \frac{\left| \frac{W' C_0 W}{n} \right|^{1/2}}{\left| \frac{W' C_1 W}{n} \right|^{1/2}} \quad (4.59)$$

2) Since σ in the two models is the same parameter, the proportionality factor is automatically the same in both cases.

3) $p(\theta_0)$ and $p(\phi_1)$ being noninformative priors for θ_0 and ϕ_1 respectively may equivalently be expressed as locally uniform priors in an approximately data translating metric (see Appendix A of

Chapter 2), i.e.

$$\begin{cases} p(\phi_0) d\phi_0 = k d\phi_0 \\ p(\phi_1) d\phi_1 = k_1 d\phi_1 \end{cases} \quad (4.60)$$

Since ϕ_0 and ϕ_1 are parameters of the same kind, we shall assume $k_0 = k_1$.

4) The locally uniform priors for the linear starting parameters M_0 and M_1 are similarly assumed to have the same constant of proportionality, as they too are parameters of the same kind

$$p(M_0|\phi_0) = p(M_1|\phi_1) \quad (4.61)$$

5) Finally prior indifference between H_0 and H_1 is expressed as

$$p(H_0) = p(H_1) \quad (4.62)$$

4.3.3 Posterior model probabilities.

In general, if a number of alternative models H_m , $m = 1, 2, \dots, N$, were under consideration, and it was known, that the data generating process was actually counted among these models, then Bayes' theorem states that

$$p(H_m|Y) = \frac{p(Y|H_m) p(H_m)}{\sum_m p(Y|H_m) p(H_m)} \quad (4.63)$$

Because such knowledge is unavailable in practice, or more fundamentally, because no model is ever exactly correct, probability statements computed from (4.63) cannot be given an absolute interpretation.

Rather they are quantitative expressions of the evidence in the data concerning the relative appropriateness of the suggested alternatives. So that even if perhaps some (unknown) model H_* might be superior to H_0 and H_1 , and in a three way comparison would have claimed most of the posterior probability, it is still legitimate to write:

$$\begin{aligned} \frac{p(H_0|Y)}{p(H_1|Y)} &= \frac{p(Y|H_0) p(H_0)}{p(Y|H_1) p(H_1)} \\ &= \frac{\int p(Y|\phi_0, H_0) p(\phi_0|H_0) p(H_0) d\phi_0}{\int p(Y|\phi_1, H_1) p(\phi_1|H_1) p(H_1) d\phi_1} \\ &= \frac{p(H_0)}{p(H_1)} \frac{\int_{M_0} I_{H_0}(\phi_0, \phi_0, \sigma) dM_0 d\phi_0}{\int_{M_1} I_{H_1}(\phi_1, \phi_1, \sigma) dM_1 d\phi_1} \\ &= \frac{\int_{\phi_0} p(\phi_0) (SS_0(\phi_0, \phi_0))^{\frac{n-p-1}{2}} d\phi_0}{\int_{\phi_1} p(\phi_1) (SS_1(\phi_1, \phi_1))^{\frac{n-p-1}{2}} d\phi_1} \quad (4.64) \end{aligned}$$

with

$$\begin{aligned} I_{H_0}(M_0, \phi_0, \sigma) &= (2\pi\sigma^2)^{-n/2} \exp(-\frac{1}{2\sigma^2} (SS_0(\phi_0, \phi_0) + (\phi_1 - M_0)^2)) |X' S_0 X|^{1/2} \sigma^{-1} p(\phi_0) \\ \text{and} \\ I_{H_1}(M_1, \phi_1, \sigma) &= (2\pi\sigma^2)^{-n/2} \exp(-\frac{1}{2\sigma^2} (SS_1(\phi_1, \phi_1) + (\phi_2 - M_1)^2 + (\phi_1 - \phi_2)^2)) |X' S_1 X|^{1/2} \sigma^{-1} p(\phi_1). \end{aligned}$$

Unlike the integrations leading to the last expression in (4.64), the still remaining integration over ϕ cannot be carried out analytically. However in terms of the data translating matrix ϕ induced by $p(\phi)$, we may write

$$\frac{p(H_0|Y)}{p(H_1|Y)} = \frac{(SS_0(\phi_0, \phi_0))^{\frac{n-p-1}{2}} \int_{\phi_0} (SS_0(\phi_0, \phi_0)/SS_0(\phi_0, \phi_0))^{\frac{n-p-1}{2}} d\phi_0}{(SS_1(\phi_1, \phi_1))^{\frac{n-p-1}{2}} \int_{\phi_1} (SS_1(\phi_1, \phi_1)/SS_1(\phi_1, \phi_1))^{\frac{n-p-1}{2}} d\phi_1} \quad (4.65)$$

and having now scaled the integrand so that its maximum value is one, it is recognized, that, by virtue of the principle on which the approximate data translating metric $\hat{\phi}$ rests (i.e. constant expected spread) the integrals in (4.65) are nearly constant, not depending on the data (see Appendix A of Chapter 2).

Hence we have established the remarkably simple

Rule: The evidence in the data, \underline{y} , concerning whether serial correlation is better accounted for by an AR-1 noise model (H_0) with a fixed mean, or by a nonstationary noise model (H_1) implying that the differenced data have AR-1 noise, is expressed through the minimum residual sum of squares of each model; specifically

$$\frac{P(H_0|\underline{y})}{P(H_1|\underline{y})} = \left(\frac{SS_0(\hat{\theta}, \hat{\phi}_0)}{SS_1(\hat{\theta}(1), \hat{\phi}_1)} \right)^{\frac{n-p-1}{2}} \quad (4.66)$$

If one had been unwilling to make the assumptions 1) through 5) in Subsection 4.3.2, then Equation (4.66) would have been derived up to a multiplicative constant κ . In favor of those assumptions is the fact, that they lead to $\kappa = 1$, an intuitively very appealing result.

4.4 The CGK data.

In this re-examination of the CGK data, we shall use all 55 observations actually listed in their paper[†].

Analyzing these data on the basis of model H_0 (4.38) the minimum residual sum of squares is $SS_0(\hat{\theta}, \hat{\phi}_0) = .80 = 15624$.

The marginal posterior distribution of $\hat{\phi}_0$, $p(\hat{\phi}_0|\underline{y})$ (4.10) is

shown in Figure 4.1a. It rather conclusively rules out the possibility[†] CGK only used the first 51 observations in their analysis, and so did BAN in their reanalysis.

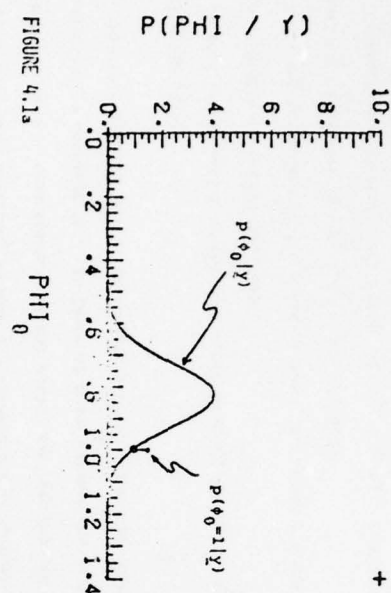
of $\hat{\phi}_0$ being zero; the body of the distribution lies in the positive stationary range of $\hat{\phi}_0$.

Posterior distributions of $\hat{\theta}_1$, $\hat{\theta}_2$ and $\hat{\theta}_3$ are given in Figures 4.2a,b and c respectively. The curves $P_0(\hat{\theta}_j|\hat{\phi}_0 = 0, \underline{y})$ show the results of relying on an assumption of independence, i.e. they correspond to the CGK analysis. The curves $P_0(\hat{\theta}_j|\hat{\phi}_0 = \hat{\phi}_0, \underline{y})$ correspond to the analysis of BAN, while the marginal distributions $P_0(\hat{\theta}_j|\underline{y})$ represent the Bayesian inference. It is seen, that the marginal posterior distributions and the posterior distribution conditional on $\hat{\phi}_0 = \hat{\phi}_0$ are almost identical for $\hat{\theta}_1$ and $\hat{\theta}_2$ and indistinguishable for $\hat{\theta}_3$. On the other hand the conditional distributions $P_0(\hat{\theta}_j|\hat{\phi}_0 = 0, \underline{y})$ differ considerably from those, both in spread and location.

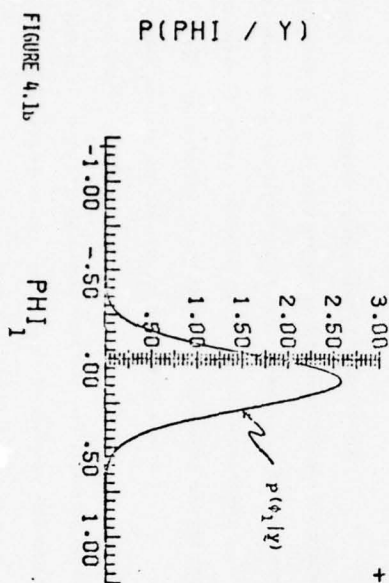
The cause of these shifts are made clear by looking at the joint posterior distribution (4.16) of $\{\hat{\theta}_j, \hat{\phi}_0\}$ $j = 1, 2, 3$ shown in Figures 4.3a,b and c. The contours of the HPD regions with probability contents 50%, 90%, 95% and 99% were computed using numerical integration, while those labeled 99.9%, 99.99%, 99.999% and 99.9999% were constructed using (4.18); and they are to be considered rough guides only since the normal approximation may be poor.

These contour plots show, that the parameters $\hat{\theta}_j$ and $\hat{\phi}_0$ are dependent. It is seen, that the distribution of $\hat{\theta}_j$ conditional on $\hat{\phi}_0$ shifts location as well as change in spread as $\hat{\phi}_0$ moves along its axis. In particular if $\hat{\phi}_0$ is set equal to the extremely unlikely value zero, then the shift is quite large. The 95% HPD intervals for $\hat{\theta}_j$, $j = 1, 2, 3$, are traced as a function of $\hat{\phi}_0$ in these contour plots Figures 4.3a,b and c (as well as in the plots of Figures 4.3d, 4.6a,

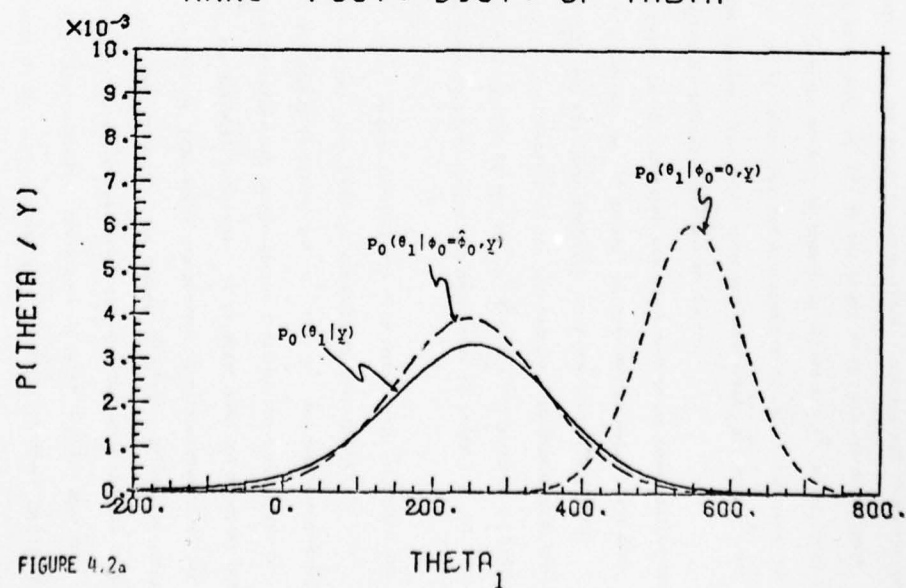
MARG. POST. DIST. OF PHI



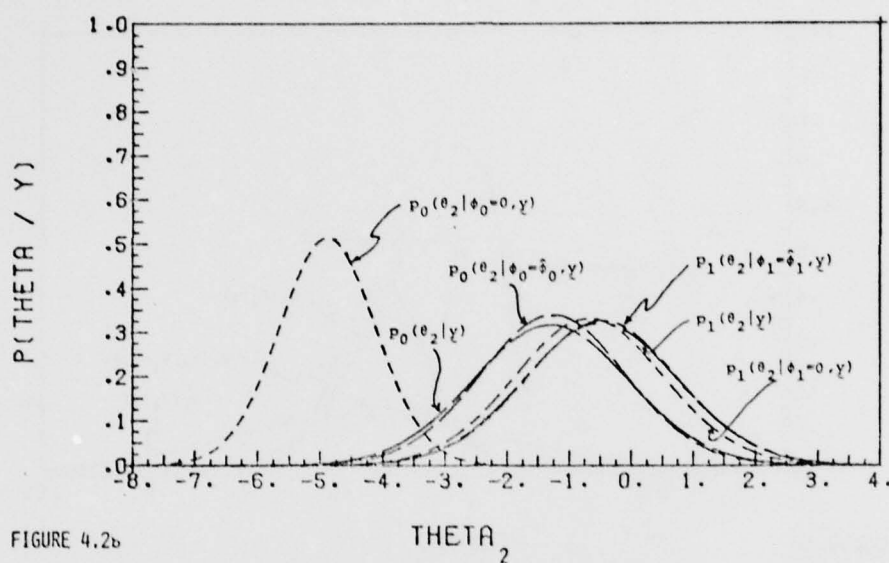
MARG. POST. DIST. OF PHI



MARG. POST. DIST. OF THETA

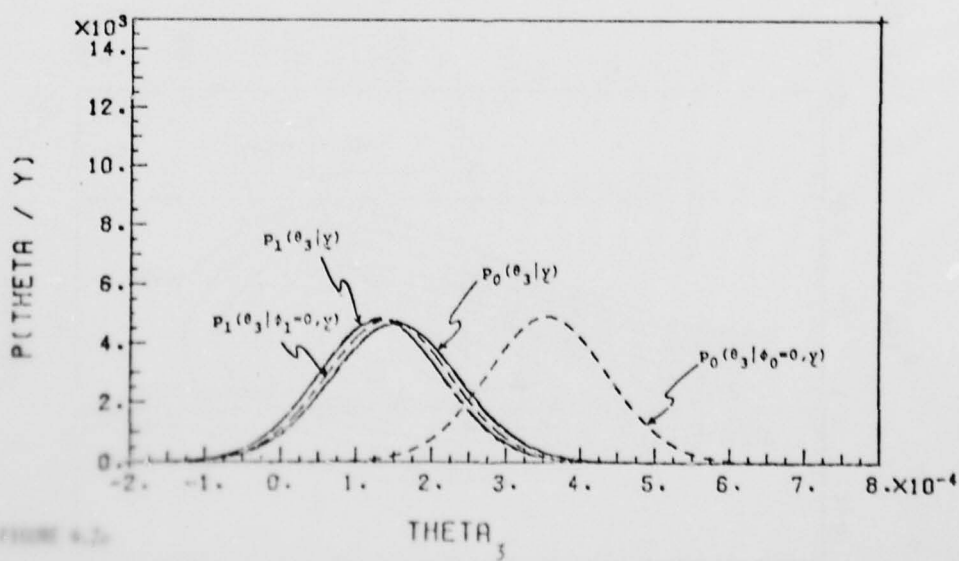


MARG. POST. DIST. OF THETA



153

MARG. POST. DIST. OF THETA



154

MARG. POST. DIST. OF THETA

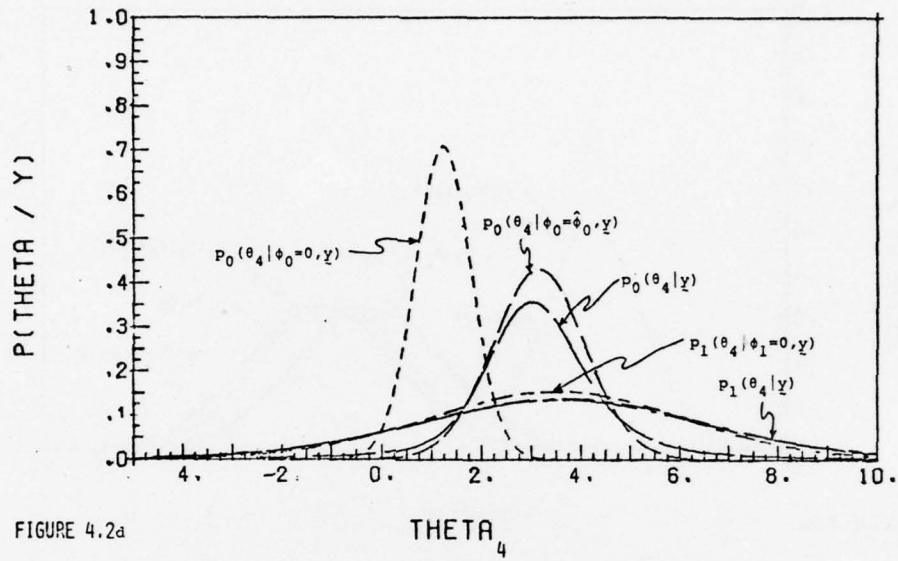


FIGURE 4.2a

THETA₄

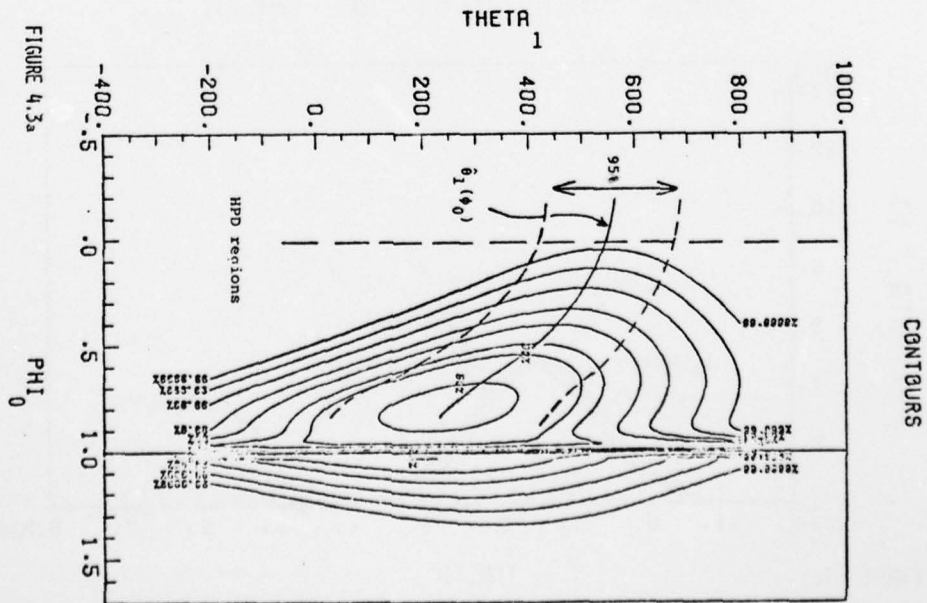


FIGURE 4.3a

CONTOURS

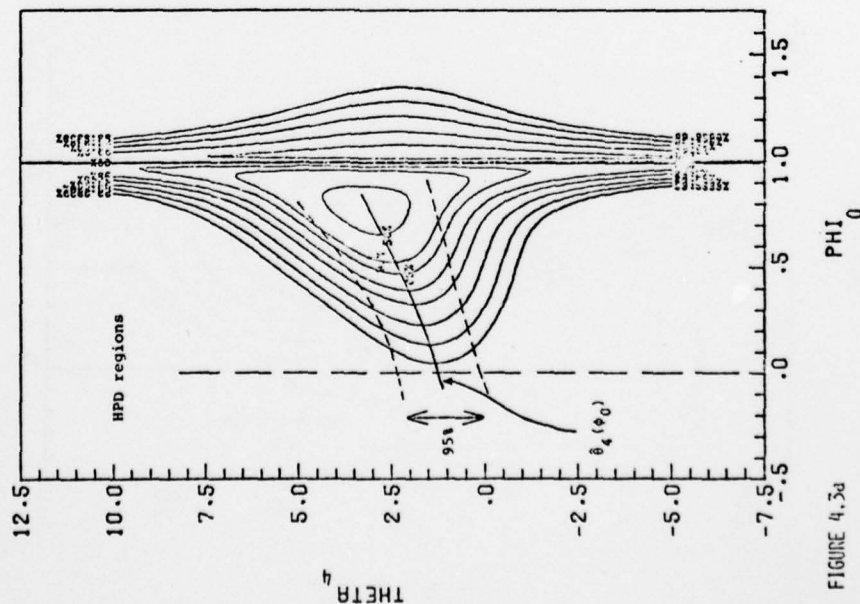


FIGURE 4.3a

b and c). It is reminded, that the 95% HPD interval for θ_j conditional on $\phi_0 = \phi_0^*$ is identical to the 95% confidence interval for θ_j which would have been derived from a sampling theory point of view assuming $\phi_0 = \phi_0^*$.

Of course integrating these joint distributions over ϕ_0 yield the marginal distributions $P_0(\theta_j)$ shown in Figures 4.2a, b and c.

As discussed in Section 4.1, the line $\phi_0 = 1$ is excluded from the (θ_j, ϕ_0) plane. Inferences about θ_2 and θ_3 conditional on $\phi_0 = 1$ is made using (4.23). The conditional posterior distribution $P_0(\theta_j | \phi_0 = 1, y)$ is found in Figures 4.2b and 4.2c as $P_1(\theta_j | \phi_0 = 1, y)$, $j = 1, 2$ (see below). As far as the mean θ_1 is concerned, this parameter vanishes when $\phi_0 = 1$. It is seen from Figure 4.3a that as ϕ_0 approaches 1, θ_1 becomes increasingly indeterminate. It is noted, that for example the 95% HPD region consists of two "islands" on either side of $\phi_0 = 1$.

To see to what extent the data supports the notion of the noise being a random walk ($\phi_0 = 1$), one may more closely examine $p(\phi_0 | y)$ Figure 4.1a at $\phi_0 = 1$. If continuity were maintained through this point, then its posterior density would be .92, but using (4.32) we find the posterior density at $\phi_0 = 1$ to be

$$p(\phi_0 = 1 | y) = .92 e^{1/2 \left(\frac{17179}{17152} \right)^2 - \frac{55-4-1}{2}} = 1.46 \quad (4.67)$$

This elevated density makes the possibility of $\phi_0 = 1$ somewhat more plausible. (It may be remarked, that to compute $\lim_{\phi_0 \rightarrow 1} SS_0(\theta, \phi_0) =$

$SS_0(\theta, \phi_0) = 17152$, the "substitution parameter" θ_1^* for the vanishing mean θ_1 , serves in this case as a deterministic trend for the

differences and not as mean for the differences, since θ_4 assume that role when $\phi_0 = 1$; see Section 3.2 of Chapter 3).

If alternatively the differenced data are analyzed along the lines of (4.34), i.e. if \tilde{y} is considered generated by H_1 (4.39), then the residual sum of squares is $SS_1(\tilde{\theta}_1, \hat{\phi}_1) = .12 = 16853$, i.e. not quite as good a fit as H_0 provided.

Figure 4.1b shows, that the autoregressive parameter $\hat{\phi}_1$ for the differences has its marginal posterior distribution located near $\phi_1 = 0$.

The marginal posterior distributions $P_1(\theta_j | \tilde{y})$, $j = 2, 3$, are shown in Figures 4.2b and 4.2c. $P_1(\theta_j | \hat{\phi}_1 = \hat{\phi}_1, \tilde{y})$, $j = 2, 3$ are also drawn, these distributions are virtually indistinguishable from the marginal distributions. Also the posterior $P_1(\theta_j | \phi_1 = 0, \tilde{y})$ are drawn, since $\hat{\phi}_1 = 0$ means that the noise \tilde{f} is a random walk, these curves may also be interpreted as $P_0(\theta_j | \phi_0 = 1, \tilde{y})$ (except one observation has been lost in the beginning of the series).

Turning to what inferences can be drawn about the trend θ_4 , we observe in Figure 4.2d, that differencing has a marked effect on what the data have to say about θ_4 . It is seen, that if model H_1 is used (or if $\phi_0 = 1$ in H_0), then the need for having a trend in the model becomes much less persuasive. The question now arises of how much the testimony of H_1 shall be discounted relative to that of H_0 , due to H_1 's slightly inferior fit. Putting the rule (4.66) to use we find

[†] These conditional distributions correspond essentially to the inference made by B&N in connection with their IMA-1 noise model, as that model is very nearly the same as the present ARIMA (1,1,0) for small ϕ_1 .

$$\frac{P(H_0 | \tilde{y})}{P(H_1 | \tilde{y})} = \left(\frac{SS_0}{SS_1} \right)^{\frac{55-4-1}{2}} = 6.64 \quad (4.68)$$

or

$$\begin{cases} P(H_0 | \tilde{y}) = .87 \\ P(H_1 | \tilde{y}) = .13 \end{cases} \quad (4.69)$$

So as far as the data are concerned it would be a mistake to completely disregard the findings associated with H_1 , in particular it may be contemplated to reanalyze the data without having a deterministic trend in the model. Before doing that we shall mention, that in situations like the present one, where the posterior model probabilities are not very conclusive, the possibility exists of combining the marginal posterior densities of the common linear parameters θ_j , $j = 2, 3, 4$ from the two models. Conditional on H_0 or H_1 being an adequate model (i.e. $P(H_0) + P(H_1) = 1$) we may write

$$P(\theta_j | \tilde{y}) = P_0(\theta_j | H_0, \tilde{y}) P(H_0 | \tilde{y}) + P_1(\theta_j | H_1, \tilde{y}) P(H_1 | \tilde{y}) \quad (4.70)$$

in other words a choice between H_0 and H_1 is not really necessary in order to make inferences about θ_j , $j = 2, 3, 4$.

Reanalyzing the data without a deterministic trend in the model, i.e. rewriting (4.1) as

$$y_i = \theta_1 + \theta_2 x_{1,i} + \theta_3 x_{2,i} + N_i \quad (4.71)$$

we find for H_0 , that the marginal posterior distribution of ϕ_0 , Figure 4.4a, has shifted dramatically towards larger values (around 1), whereas $P(\phi_1 | \tilde{y})$ of model H_1 has hardly changed at all, Figure 4.4b.

Looking at the density $P(\phi_0 | \tilde{y})$ at $\phi_0 = 1$, we see that if continuity were maintained then this point would have a density of about 6.9, but using (4.32) we find

$$p(\hat{\phi}_0=1|\tilde{y}) = 6.9 e^{1/2 \left(\frac{17895}{17179} \right)} - \frac{55-3-1}{2} = 4.0. \quad (4.72)$$
 Hence the inducement to difference has increased by the elimination of the trend from the model, but the increase is not as pronounced as it would have appeared, if $p(\hat{\phi}_0|\tilde{y})$ had been interpreted as going continuously through $\hat{\phi}_0=1$.

Applying the rule (4.66) we also find, that upon elimination of θ_4 , the model structure H_1 gains support vis-a-vis H_0 . Specifically

$$\frac{P(H_0|\tilde{y})}{P(H_1|\tilde{y})} = \left(\frac{SS_0(\hat{\theta}_0, \hat{\phi}_0 = .96)}{SS_1(\hat{\theta}_1, \hat{\phi}_1 = .15)} \right)^{\frac{55-3-1}{2}} = \left(\frac{16911}{17379} \right)^{\frac{51}{2}} = 2.00 \quad (4.73)$$

or

$$\begin{cases} P(H_0|\tilde{y}) = .67 \\ P(H_1|\tilde{y}) = .33 \end{cases} \quad (4.74)$$

For illustration Figures 4.6a, b and c show the joint posterior distribution of $\{\theta_j, \hat{\phi}_0\}$ $j = 1, 2, 3$ in relation to model H_0 . Now as before these contour plots demonstrate the dependence of θ_j on $\hat{\phi}_0$. Integrating the joint distributions over $\hat{\phi}_0$ yields the marginal posteriors $p_0(\theta_j|\tilde{y})$ $j = 1, 2, 3$ drawn in Figures 4.5a, b and c. These figures also show the posterior distributions conditional on $\hat{\phi}_0=0$ and $\hat{\phi}_0=\hat{\phi}_0$, as well as the marginal and conditional distributions of θ_j , $j = 2, 3$ relating to model H_1 . It is noted, that whatever disagreement existed between $p_0(\theta_j|\tilde{y})$ vs. $p_1(\theta_j|\tilde{y})$, $j = 2, 3$, has now disappeared.

The marginal posterior distributions of θ_2 and θ_3 concur with the conclusion drawn by B&W, that "there is no real evidence of any

MARG. POST. DIST. OF PHI

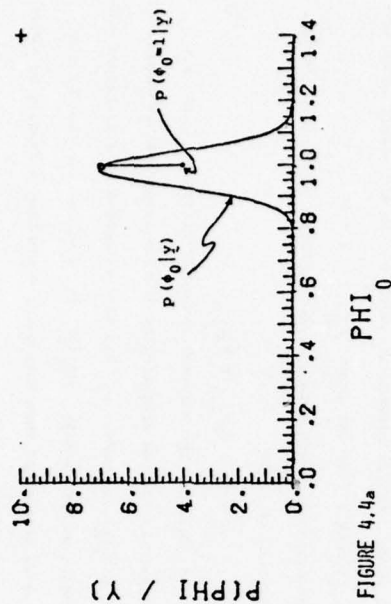


FIGURE 4.4a

MARG. POST. DIST. OF PHI

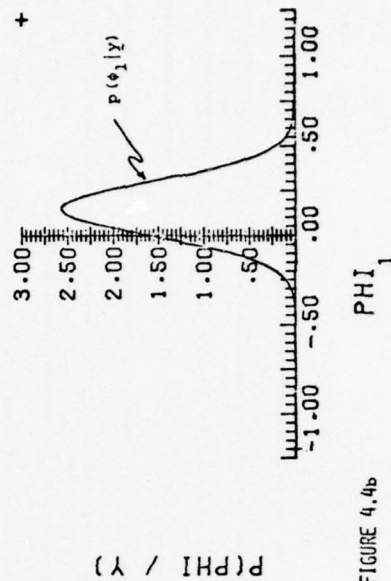


FIGURE 4.4b

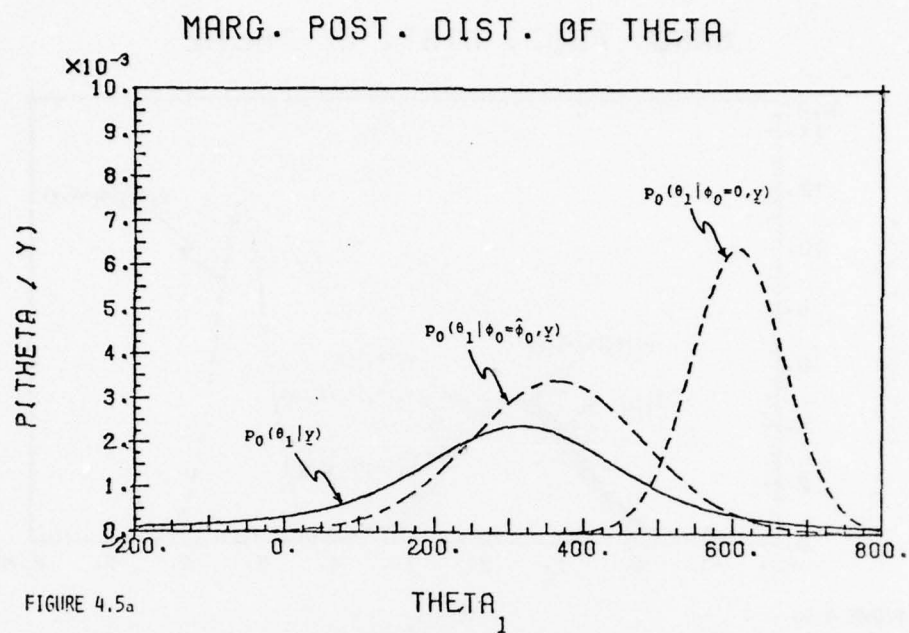


FIGURE 4.5a

165

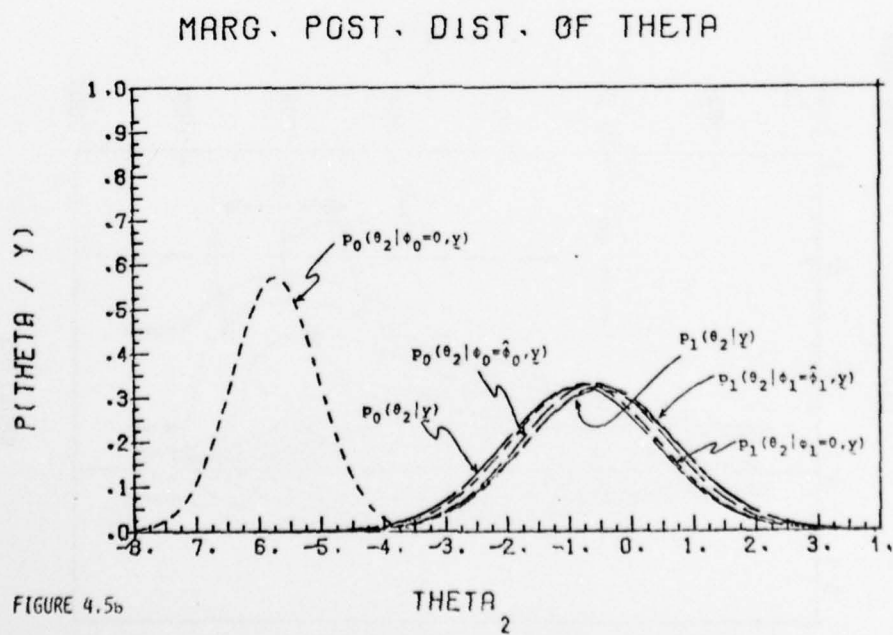
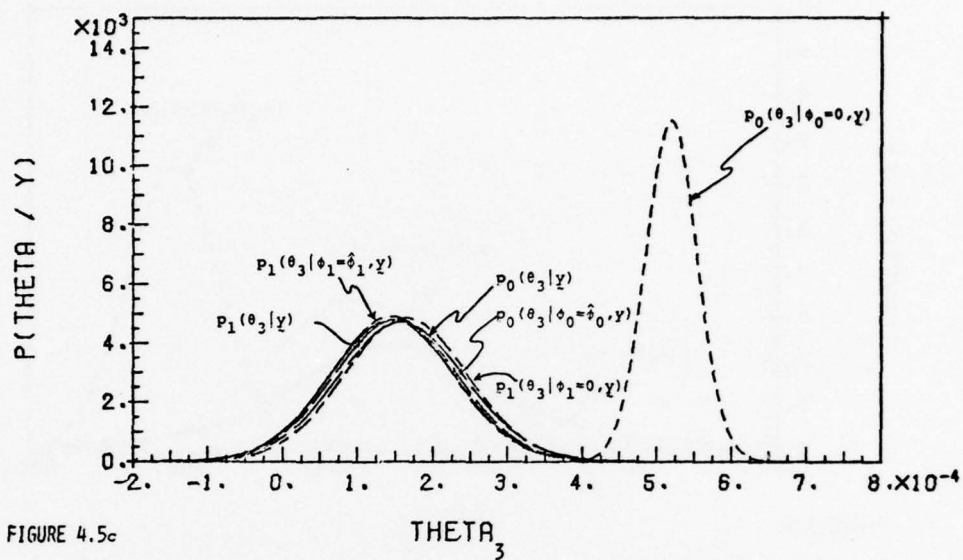


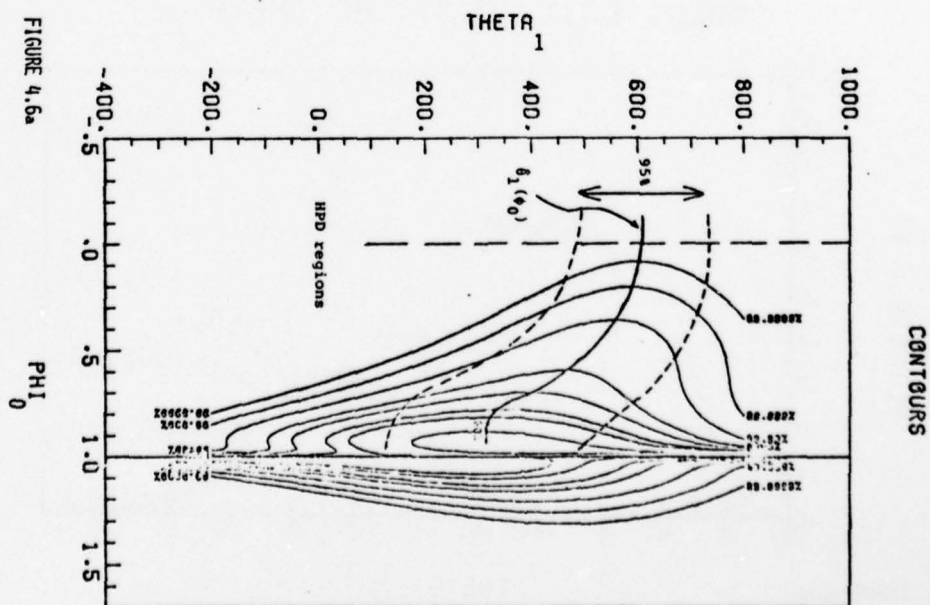
FIGURE 4.5b

166

MARG. POST. DIST. OF THETA

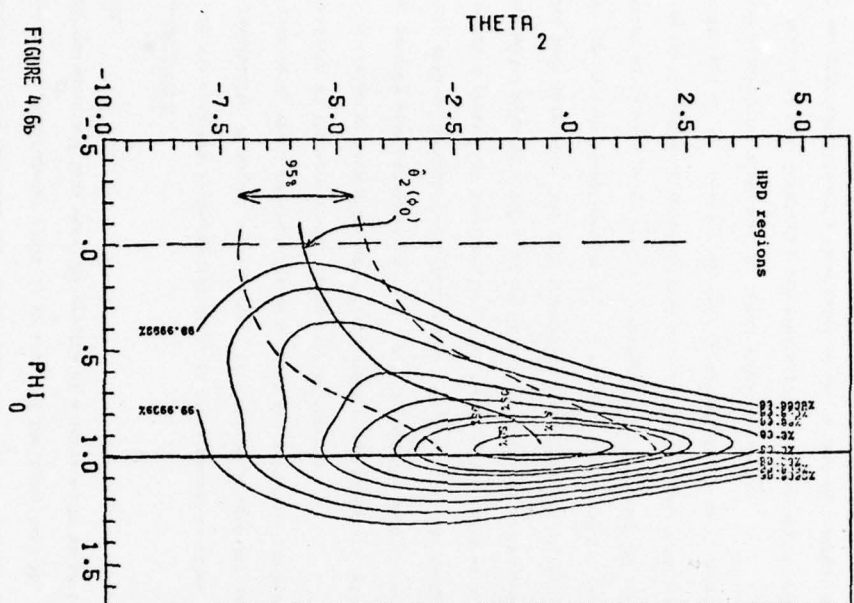


167

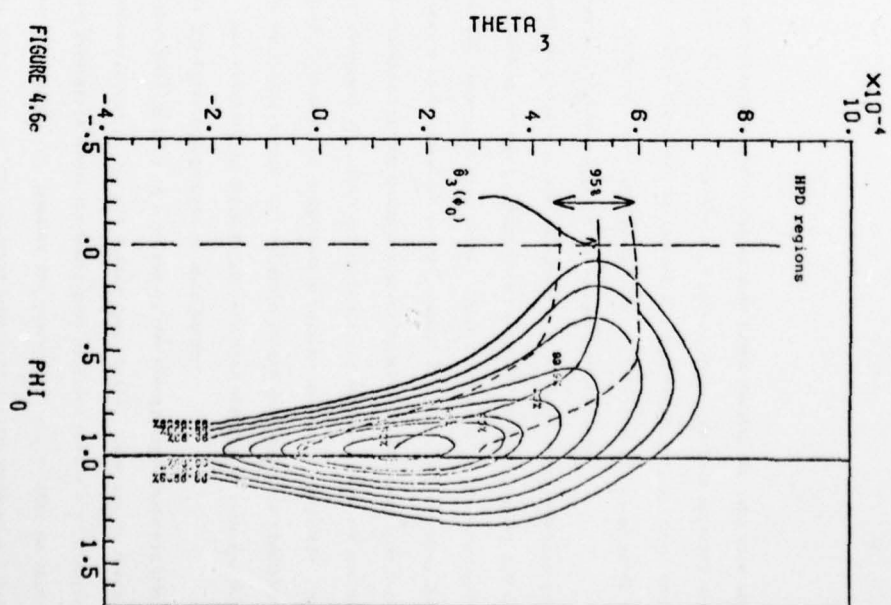


168

CONTOURS



CONTOURS



Of course the AR-1 process includes a random walk as a special case, namely for $\phi = 1$. To assess the plausibility of that particular choice for ϕ the marginal posterior distribution of ϕ , $p(\phi|y)$, should be looked at. However the density at $\phi = 1$ must be determined separately since the model mean vanishes when $\phi = 1$, creating a singularity point. It was argued that $p(\phi|y)$ does not in general go continuously through $\phi = 1$, and it was specifically proposed how to find the posterior density in that point.

The question of differencing was also approached from the point of view of seeing whether the evidence in the data supports spending one degree of freedom on differencing rather than on a fixed mean. In this treatment the ratio of the posterior probabilities for two reasonable alternative one parameter noise models was derived. One model assumes a fixed mean for an AR-1 noise, the other implies that the differenced series has AR-1 noise. In situations of doubt about which model should be the preferred one, marginal inference about the common linear parameters in the two models is actually still possible in the Bayesian framework without having to make a choice.

It would seem, that giving the data an opportunity to speak out in this quantitative way on whether a differencing is justified, would extend the appropriateness of a linear model analysis allowing for serial correlation to cover an even wider variety of real data sets.

relation between the output on the one hand and the two lagged input on the other". The tail area for $p(\theta_2|y)$ to the right of $\theta_2 = 0$ is about 28%, and the tail area to the left of $\theta_3 = 0$ is about 4%, which is not particularly striking in light of the fact, that the proposed variables and lags were the results of a large scale selection effort.

4.5 Conclusion.

In the analysis of sequential data, it is necessary to allow for the possibility of serial correlation. For example in economic and business data, where observations have to be taken in time order autocorrelation of the errors will ordinarily be expected.

Generally an assumption of a first order autoregressive (AR-1) noise process with parameter ϕ ($-\infty < \phi < \infty$) will be much closer to reality than an assumption of white noise. Analyzing in the Bayesian framework a particular data set (the CGK data) on the basis of a linear model with AR-1 noise (along the lines developed in Chapter 2), it was made very clear how very misleading inferences may be reached about the regression parameters θ_j , $j = 1, 2, \dots, p$, if serial correlation is ignored (i.e. $\phi = 0$). Specifically it was seen, by looking at the joint posterior distributions of $\{\theta_j, \phi\}$ that the conditional distribution of θ_j changes not only in spread, but also in location as the conditional value for ϕ moves along its axis.

Realizing that dependence can create problems, one way in which data analysts traditionally have tried to get around this complication, is by differencing and then assuming, that the first differences have white noise errors, i.e. that the errors in the original observations follow a random walk.

For this example it is found that

$$P_n(\phi = .5|\tilde{y}) = 1.06 e^{1/2} .991 = 1.76 \quad (4A.4)$$

so that the point $\phi = .5$, which is the true value for ϕ , actually has the highest posterior density of any point.

Appendix, Other applications

In the analysis of the spirits data in Section 3.5.3 of Chapter 3 it might be wondered whether or not these data should be differenced, since the marginal posterior distribution of the autoregressive parameter ϕ , labeled $P_n(\phi|\tilde{y})$ in Figure 3.6b, is located near $\phi = 1$. The density $\lim_{\phi \rightarrow 1} P_n(\phi|\tilde{y})$ is 8.92. Using (4.32) we find that the adjusted density at $\phi = 1$ is

$$p(\phi=1|\tilde{y}) = 8.92 e^{1/2} \left(\frac{-.03298}{.03009} \right)^2 = .74. \quad (4A.1)$$

This much lower density discourages a differencing.

If the spirits data are analyzed according to model H_1 , then the residual sum of squares is $SS_1 = .03284$, so that using (4.66) we find

$$\frac{P(H_0|\tilde{y})}{P(H_1|\tilde{y})} = \left(\frac{-.02972}{.03284} \right)^2 = \frac{26}{1}, \quad (4A.2)$$

which also testifies strongly against differencing.

An interesting application of the continuity adjustment derived in Section 4.2 presents itself in connection with the "artificial data" analyzed in Section 2.6.1 of Chapter 2. Figure 2.5b shows the marginal posterior distribution $P_n(\phi|\tilde{y})$; however since $\phi = .5$ is a singularity point, the density at that point must be found separately, namely from

$$p(\phi=.5|\tilde{y}) = \lim_{\phi \rightarrow .5} p(\phi|\tilde{y}) = \frac{1}{2} \left(\frac{SS(\theta_{(1)}, \phi=.5)}{SS(\theta, \phi=.5)} \right)^{\frac{n-p-1}{2}} \quad (4A.3)$$

where $SS(\theta_{(1)}, \phi=.5)$ and $SS(\theta, \phi=.5) = \lim_{\phi \rightarrow .5} SS(\theta, \phi)$ are computed as explained in Section 3.2 of Chapter 3.

As far as $\hat{\theta}$ is concerned, it seems to make little difference in the analysis of real data however, if instead a convenient uniform prior for $\hat{\theta}$ is adopted. The integration over $\hat{\theta}$ leading to the marginal posterior distribution of $\hat{\theta}$, must be carried out numerically, but excellent approximations may be obtained from a weighted sum of, say, five t-distributions.

Recognizing that the assumption of independence can be a crucial one, it has become a widespread practice in regression work, when observations are made in time order, to test for serial correlation using the well known Durbin Watson statistic. This test corresponds to determining whether a certain estimate, $\hat{\rho}$ of ρ is significantly different from zero. Using a likelihood or Bayesian approach it is seen in Chapter 3 how better estimates for ρ are obtained, and how approximate confidence or HPD intervals for ρ are constructed. These approaches suggest testing procedures for the general hypothesis $\rho = \rho_0$; and in particular for $\rho_0 = 0$ the alternatives may be viewed as competitors to the DW-test. It is demonstrated by simulation, that unless the linear model lacks a mean, the DW-test is comparable in power to the suggested alternatives. It is argued however, that it is really beside the point to test a null hypothesis of independence in such situations where serial correlation is to be expected. It was illustrated, that proceeding in the fashion implied by the DW testing approach carries a penalty, as namely the further analysis concerning $\hat{\theta}$ may give very misleading results both if inferences about $\hat{\theta}$ are made conditionally on $\rho = 0$ (when that hypothesis is "accepted") or if they are drawn conditionally on $\rho = \hat{\rho}$ (if $\rho = 0$ is "rejected").

CHAPTER 5

Summary

The analysis of linear models with independent, homoscedastic, normal noise (white noise) occupies a prominent position in statistics, and is used extensively as a tool in applied work in fields like business, economics and sociology as well as in the natural sciences and engineering. However when data have been collected sequentially in time or space the assumption of independence will ordinarily be unrealistic. If the assumptions are relaxed by introducing a single parameter to allow for serial correlation, then many real life data sets with autocorrelated errors may be approximately modelled.

Assuming that the observations (original or transformed) can adequately be represented by a linear model (with parameters $\hat{\theta}$) whose noise term follows a first order autoregressive scheme (with parameter ρ), it was shown in Chapter 2 how in the Bayesian framework inferences can be drawn about $\hat{\theta}$ and ρ jointly, conditionally and marginally. Two AR-1 structures were considered; one (employing an additional starting parameter) covers explosive as well as stationary cases, the other assumes stationarity and reversibility. Approximately noninformative prior distributions for the parameters complementing the likelihood functions are suggested, in particular it is argued, that prior independence between $\hat{\theta}$ and ρ is not an appropriate assumption. In situations where the noise is clearly stationary, the inference about $\hat{\theta}$ and in turn about ρ appears very little affected by which of the two AR-1 schemes is employed. Using the more flexible not necessarily stationary model is generally recommended, although special care must be exercised in applying the locally approximately noninformative prior for ρ .

The Bayesian analysis of the regression model with AR-1 noise

solves these problems. As far as $\hat{\theta}$ is concerned, it is proposed that the full Bayesian analysis may be approximated by a conditional analysis using the maximum likelihood estimate $\hat{\phi}$ of ϕ . ($\hat{\phi}$ is found by estimating ϕ and $\hat{\theta}$ simultaneously by least squares.) The resulting conditional inference is paralleled from a sampling theory point of view when inferences about $\hat{\theta}$ are drawn as if $\phi = \hat{\phi}$.

The difficulties referred to above arise because of dependence of $\hat{\phi}$ and $\hat{\theta}$. The nature of this dependence is further clarified in Chapter 4, where for a particular example the joint posterior distributions of (θ_j, ϕ) $j = 1, 2, \dots, p$ were considered. These plots demonstrate that the posterior distribution of θ_j conditional on $\phi = \phi^0$ not only changes in spread as ϕ^0 moves along the ϕ -axis, but also shifts location. One way in which data analysts traditionally have tried to get around the problems of serial correlation is by differencing and then assuming that the first differences have white noise errors, i.e. that the errors in the original observations follow a random walk. This is equivalent to assuming that $\phi = 1$ in the AR-1 noise model, and the plausibility of that particular value for ϕ may be assessed by studying the marginal posterior distribution of ϕ . For models involving a mean this parameter vanishes for $\phi = 1$ creating a singularity point; and it is proposed how the density may be determined in such (distinct) points. The question of differencing was also approached in Chapter 4 from the point of view of determining whether the evidence in the data favors expending one degree of freedom on differencing or on a fixed mean. Specifically a quantitative rule was derived which expresses (in terms of relative posterior model probabilities) the evidence for and

against differencing prior to an analysis of the data along the lines of Chapter 2.

REFERENCES

1. Abraham, B. and G. E. P. Box (1975), "Linear Models, Time Series and Outliers, 2: Outliers in Linear Models", Technical Report No. 437, Dept. of Statistics, Univ. of Wis., Madison, Wis.
2. Abraham, A. P. J. and A. S. Louter (1971), "On a New Test for Autocorrelation in Least Squares Regression", *Biometrika*, 58, 53-60.
3. Aitken, A. C. (1935), "On least squares and linear combinations of observations", *Proc. Roy. Soc. Edinburgh*, 55, 42-48.
4. Anderson, O. D. (1975), "A note on differencing autoregressive moving average (p,q) processes", *Biometrika*, 62, 521-523.
5. Anderson, R. L. (1954), "The Problem of Autocorrelation in Regression Analysis", *JASA*, 49, 113-129.
6. Anscombe, F. J. (1960), "Rejection of outliers", *Technometrics*, 2, 123-146.
7. Anscombe, F. J. and Tukey, J. W. (1963), "The Examination and Analysis of Residuals", *Technometrics*, 5, 141-160.
8. Bartlett, M. S. (1947), "The Use of Transformations", *Biometrics*, 3, 39-52.
9. Berenblut, I. I. and G. I. Webb (1973), "A New Test for Autocorrelated Errors in the Linear Regression Model", *JRSS series B*, 35, 33-50.
10. Box, G. E. P. (1976), personal communication.
11. Box, G. E. P. and D. R. Cox (1964), "An Analysis of Transformations", *JRSS Series B*, 26, 211-252.
12. Box, G. E. P. and Henson, T. L. (1969), "Model fitting and discrimination", Tech. Report No. 211, Dept. of Stat., Univ. of Wis., Madison, Wis.
13. Box, G. E. P. and Hill, W. J. (1967), "Discrimination among mechanistic models", *Technometrics*, 9, 57-71.
14. Box, G. E. P. and G. M. Jenkins (1970), "Time Series Analysis, Forecasting and Control", Holden-Day.
15. Box, G. E. P. and P. Newbold (1971), "Some Comments on a Paper of Coen, Gome and Kendall", *JRSS series A*, 134, 229-240.
16. Box, G. E. P. and G. C. Tiao (1973), "Bayesian Inference in Statistical Analysis", Addison-Wesley Publishing Company.
17. Champenowne, D. G. (1948), "Sampling Theory Applied to Autoregressive Sequences", *JRSS series B*, 10, 204-242.
18. Cochran, D. and G. H. Orcutt (1949), "Application of Least Squares Regression to Relationships Containing Auto-Correlated Error Terms", *JASA*, 44, 32-61.
19. Coen, P. J., E. D. Gome and M. G. Kendall (1969), "Lagged Relationships in Economic Forecasting", *JRSS series A*, 132, 133-152.
20. Durbin, J. (1969), "Testing for serial correlation in regression analysis based on the periodogram of least squares residuals", *Biometrika*, 56, 1-15.
21. Durbin, J. (1970), "An Alternative to the Bounds Test for Testing for Serial Correlation in Least Squares Regression", *Econometrica*, 38, 427-429.
22. Durbin, J. (1970a), "Testing for Serial Correlation in Least-Squares Regression when some of the Regressors are Lagged Independent Variables", *Econometrica*, 38, 410-421.
23. Durbin, J. and G. S. Watson (1950), "Testing for Serial Correlation in Least Squares Regression. I.", *Biometrika*, 37, 409-428.

AD-A054 556

WISCONSIN UNIV-MADISON MATHEMATICS RESEARCH CENTER
STUDIES IN THE ANALYSIS OF SERIALY DEPENDENT DATA.(U)
MAR 78 L C PALLESEN

F/G 12/1

DAA629-75-C-0024

NL

UNCLASSIFIED

MRC-TSR-1837

2 OF 2
AD
A054556

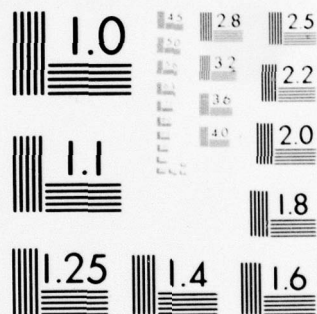


END

DATE
FILMED

7 - 78

DDC



24. Durbin, J. and G. S. Watson (1951), "Testing for Serial Correlation in Least Squares Regression. II.", *Biometrika*, 38, 159-178.
25. Durbin, J. and G. S. Watson (1971), "Testing for Serial Correlation in Least Squares Regression, III.", *Biometrika*, 58, 1-19.
26. Durbin, J. (1960), "Estimation of Parameters in Time-Series Regression Models", *JRSS series B*, 22, 139-153.
27. Fraser, D. A. S. (1968), "The Structure of Inference", John Wiley and Sons, Inc.
28. Grady, R. C. (1970), "Relative Efficiency of Count of Sign Changes for Assessing Residual Autocorrelation in Least Squares Regression", *Biometrika*, 57, 123-127.
29. Hannan, E. J. (1957), "Testing for Serial Correlation in Least Squares Regression", *Biometrika*, 54, 57-66.
30. Haq, M. S. (1970), "Structural analysis for the first order autoregressive stochastic process models", *Ann. Math. Statist.*, 41, 970-978.
31. Haq, M. S. (1971), "Structural inference for the linear model with autoregressive error", *J. Stat. Res.*, 5, 10-22.
32. Hildreth, C. (1969), "Asymptotic Distribution of Maximum Likelihood Estimators in a Linear Model with Autoregressive Disturbances", *The Annals of Math. Stat.*, 40, 583-594.
33. Jeffreys, H. (1961), "Theory of Probability" 3rd ed., Oxford University Press.
34. Kanemasu, H. (1973), "Topics in Model Building", Ph. D. thesis, Dept. of Stat., Univ. of Wis., Madison, Wis.
35. Levenbach, H. (1972), "Estimation of autoregressive parameters from a marginal likelihood function", *Biometrika*, 59, 61-71.

36. Orcutt, G. H. and Cochrane, D. (1949), "A Sampling Study of the Merits of Autoregressive and Reduced Form Transformations in Regression Analysis", *JASA*, 44, 356-372.
37. Pierce, D. A. (1971), "Least squares estimation in the regression model with autoregressive-moving average errors", *Biometrika*, 58, 299-321.
38. Rao, C. R. (1965), "Linear Statistical Inference and its Applications", John Wiley and Sons, Inc.
39. Reiser, B. (1975), "Structural Inference for Linear Regression with Autocorrelated Errors", *Stat. Hefte*, 16, 85-104.
40. Schmidt, P. (1972), "A Generalization of the Durbin-Watson Test", *Australian Economic Papers*, 11, 203-209.
41. Smith, V. K. (1976), "The Estimated Power of Several Tests for Autocorrelation with Non-First-Order Alternative", *JASA*, 71, 879-883.
42. Sredni, J. (1970), "Problems of Design, Estimation, and Lack of Fit in Model Building", Ph.D. thesis, Dept. of Stat., Univ. of Wis., Madison, Wis.
43. Student (Cosset, W. S.) (1914), "The Elimination of Spurious Correlation Due to Position in Time or Space", *Biometrika*, 10, 179-181.
44. Theil, H. and A. L. Nagar (1961), "Testing the Independence of Regression Disturbances", *JASA*, 56, 793-806.
45. Tukey, J. W. (1957), "On the Comparative Anatomy of Transformations", *Annals Math. Stat.*, 28, 602-632.
46. Tukey, J. W. (1960), "A survey of sampling from contaminated distributions", *Contributions to Probability and Statistics: Volume Dedicated to Harold Hotelling*. Stanford University Press.

47. Watson, G. S. (1967), "Linear Least Squares Regression", Ann. Math. Stat., 38, 1679-1699.
48. Whittle, P. (1953), "Estimation and Information in Stationary Time Series", Arkiv för Matematik, 2, 423.
49. Wilks, S. S. (1962), "Mathematical Statistics", John Wiley and Sons, Inc.
50. Zellner, A. and G. C. Tiao (1964), "Bayesian Analysis of the Regression Model with Autocorrelated Errors", JASA, 59, 763-768.

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)		REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM	
1. REPORT NUMBER	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER			
1837 V					
4. TITLE (and Subtitle)		5. TYPE OF REPORT & PERIOD COVERED			
STUDIES IN THE ANALYSIS OF SERIALLY DEPENDENT DATA		Summary Report - no specific reporting period			
6. AUTHOR(s)		7. PERFORMING ORG. REPORT NUMBER			
Lars Chr. Pallesen		DAAG29-75-C-0024 ✓			
8. PERFORMING ORGANIZATION NAME AND ADDRESS		9. PROGRAM ELEMENT PROJECT, TASK AREA & WORK UNIT NUMBERS			
Mathematics Research Center, University of Wisconsin 610 Walnut Street Madison, Wisconsin 53706		#4 - Probability, Statistics and Combinatorics			
10. CONTROLLING OFFICE NAME AND ADDRESS		11. REPORT DATE			
U. S. Army Research Office P. O. Box 12211 Research Triangle Park, North Carolina 27709		March 1978			
12. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		13. NUMBER OF PAGES			
		183			
		14. SECURITY CLASS. (of this report)			
		UNCLASSIFIED			
15. DISTRIBUTION STATEMENT (of this Report)		16. DECLASSIFICATION/DOWNGRADING SCHEDULE			
Approved for public release; distribution unlimited.					
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)					
18. SUPPLEMENTARY NOTES					
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)		Regression, Linear models, Bayesian approach, Serial correlation, Autoregressive noise, Nonstationary noise, differencing			
20. ABSTRACT (Continue on reverse side if necessary and identify by block number)		<p>The analysis of linear models with independent homoscedastic, normal noise (white noise) occupies a prominent position in applied statistics. This report is concerned with the linear model analysis of data which cannot be assumed statistically independent because the data have been collected sequentially in time or space.</p> <p>Assuming that the noise in a linear model (with p-dimensional parameter vector θ) follows a first order autoregressive scheme (with parameter ϕ) it is developed in Chapter 2 how in the Bayesian framework inferences can be drawn about</p>			
DD FORM 1 JAN 73 1473 EDITION OF 1 NOV 68 IS OBSOLETE		UNCLASSIFIED			
SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)					

Abstract (continued)

ρ and ϕ jointly, conditionally and marginally. Two AR-1 schemes are considered, one covering explosive as well as stationary situations, the other assumes stationarity and reversibility a priori. The important task of choosing an appropriate joint prior distribution for the parameters is given special attention.

Recognizing that the assumption of independence can be a crucial one, it has become a widespread practice in much regression work, where observations are made in time order, to test for serial correlation using the well known Durbin-Watson test. This test corresponds to determining whether a certain estimate of ρ is significantly different from zero. It is shown in Chapter 3 that only in relation to a model lacking a mean do suggested alternative testing procedures show decisively greater empirical power than the DW-test. However it is argued that tests of a null hypothesis of independence should really not be carried out when serial correlation is to be expected. It is demonstrated that inferences about ρ may be quite misleading if made conditionally on $\phi = 0$ (when that hypothesis is "accepted") or conditionally on ϕ (when it is "rejected"). The Bayesian analysis does not suffer from these handicaps. The Bayesian inference about ρ may be approximated by a conditional inference using the maximum likelihood estimate $\hat{\phi}$ of ϕ ($\hat{\phi}$ is found by estimating ϕ and θ simultaneously by least squares), and this conditional analysis is paralleled in sampling theory framework when inferences about θ are drawn as if $\phi = \hat{\phi}$.

One way in which data analysts traditionally have tried to get around the problem of dependence, is by differencing and then assuming that the errors of the differences are independent. This is equivalent to assuming that $\phi = 1$ in the AR-1 noise model; and the plausibility of this particular value for ϕ may be assessed by studying the marginal posterior distribution of ϕ . For models involving a mean this parameter vanishes for $\phi = 1$. This creates a singularity and it is shown in Chapter 4, how the density may be found at such (distinct) points. It is also developed in the Bayesian framework, how the appropriateness of differencing may be expressed as posterior model probabilities for two alternative noise models, one assuming the original observations to have AR-1 noise, the other that the differences have AR-1 noise.